

# Conserving Memory Bandwidth in Chip Multiprocessors with Runahead Execution

Martin Karlsson and Erik Hagersten  
Department of Information Technology, Uppsala University  
P.O. Box 337, SE-751 05, Uppsala, Sweden  
{martin.karlsson}@it.uu.se

## ABSTRACT

The introduction of chip multiprocessors (CMPs) presents new challenges and trade-offs to computer architects. Architects must now strike a balance between the number of cores per chip versus the amount of on-chip cache and available pin bandwidth. Technology projections predict that the cost of pin bandwidth will increase significantly and may therefore limit the number of processor cores per CMP.

We observe a trend in many processor designs towards larger cache blocks for the highest level on-chip cache. A large cache block size is beneficial for workloads with a high amount of spatial locality. Our study confirms previous observations finding that significant parts of medium-sized cache blocks that are brought on-chip often remain unused and therefore wastefully consume pin bandwidth, especially for the commercial workloads studied. In this paper we target this waste by proposing a method of fine-grained fetches.

In this paper we show that due to characteristics of runahead execution it is possible to remove the implicit assumption that programs exhibit abundant spatial locality, with a limited performance impact. We demonstrate, using execution-driven full system simulation, that our method of fine-grained fetching can obtain significant performance speedups in bandwidth constrained systems but also yield stable performance systems that are not bandwidth limited.

## 1. INTRODUCTION

The continued decrease in transistor size combined with the increasing delay of wires relative to transistor switching speeds has led to the development of chip multiprocessors (CMPs). As more and more transistors become available per chip, it is possible to fit more and more cores per die. However, the number of off-chip signal pins<sup>1</sup> is not growing

<sup>1</sup>Throughout this paper we use the term signal pins, although signal pads may be more accurate given current packaging technology

Technical report 2005-040  
Department of Information Technology  
Uppsala University

nearly as fast leading to an increasing disparity between the number of cores that fit on a die and the available chip bandwidth. While chip and memory bandwidth since long have been a performance-sensitive resource, it has now surfaced as one of the major performance limiters.

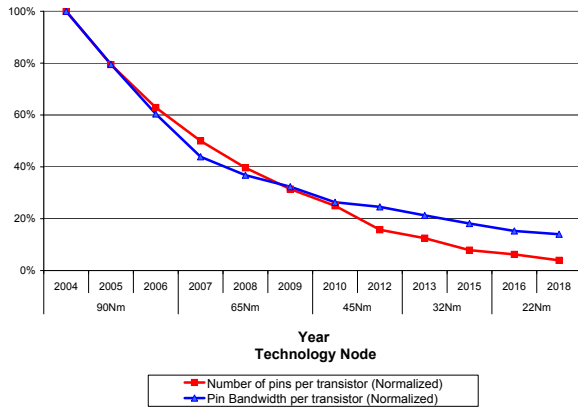
Large-scale chip multiprocessors are desirable due to their low intra chip communication cost, enabled by shared higher level caches. For many commercial applications with an abundance of thread-level parallelism[22], adding cores or threads yields a very appealing performance improvement per mm<sup>2</sup>. We therefore believe that the number of cores per chip will be scaled up as high as the chip power and bandwidth will support. In this paper we focus on conserving bandwidth in order to reduce cost and allow more processor cores per CMP. Increasing the on-chip cache reduces the bandwidth consumption. However, the performance effect of adding more cache is non-linear and for certain applications does not improve performance at all. Hence increasing on-chip cache can alleviate but not solve the bandwidth problem.

As the memory wall problem has come to overshadow other aspects of processing, various forms of run-ahead execution have been proposed[18][11][7][3][4]. Runahead execution attempts to reduce the effect of the long memory latencies by increasing the memory-level parallelism. The strength of runahead execution is its ability to prefetch data far ahead of the stall point. We will show in this paper how characteristics of runahead execution can be exploited to reduce bandwidth consumption.

## 2. TECHNOLOGY TRENDS

The off-chip bandwidth of a chip is determined by the number of signal pins and the off-chip frequency. Reliability, power and especially cost restrict the number of pins per chip[5].

The International Technology Roadmap for Semiconductors published by the Semiconductor Industry Association (SIA) contains yearly predictions of the cost-efficient number of transistors and signal pins per chip for microprocessor chip designs [21]. The 2003 edition reported an annual growth in the maximum number of signal pins that is deemed cost-effective per chip of 5% (Avg.) between 2003 and 2018. In the 2004 update, however, the projection now predicts zero growth during the same time span. This indicates a significant challenge ahead to increase the chip bandwidth



**Figure 1: The Semiconductor Industry Association’s prediction of bandwidth per transistor.**

in a cost efficient way.

SIA predicts the number of transistors per chip to grow annually by 26% (Avg.) during the same period. Hence the number of signal pins per transistor will decrease. As can be seen in Figure 1 the transistor count per signal pin ratio is projected to increase 20 fold to the year 2018. When taking relative on/off chip frequency development into account the increase is reduced to a factor of 7. This is due to the fact that signal pin frequencies are predicted to grow faster than on-chip frequencies. Several promising approaches, including optics and proximity communication, lie on the pin technology horizon. However in both cases we believe that such paradigm shifts may provide a one-time bandwidth increase, not a fundamental change in the bandwidth growth rate. We therefore believe that in the long-term perspective, the off-chip bandwidth will decrease in relation to the number of transistors per chip. If the number of processor cores and cache per chip is going to be scaled linearly with the transistors made available by future process generations, it will lead to bandwidth requirements that will be hard to satisfy.

### 3. A TREND OF GROWING CACHE BLOCKS

The lack of growth in off-chip bandwidth will force architects to reevaluate memory system design from a bandwidth-centric view. This contradicts another trend we have observed in processor design. A trend towards larger and larger cache block sizes in the highest level on-chip caches. Examples of this is observed in the Pentium, Itanium and SPARC64 processor families. The L2 cache block size increased from 32 byte to 128 between the Pentium III and 4 processors. Similarly in Itanium it increased from 64 bytes per on-chip L3 block to 128 byte per block between Merced and McKinley. The largest increase was observed in SPARC64 between version V and VI, which went from 64 bytes to 256 bytes per cache block. These increased L2 cache block sizes were undoubtedly chosen to amortize the cost of longer and longer memory latencies. However, unless all of the fetched data are used, larger cache blocks leads to some data being unnecessarily brought on chip consuming

the scarce bandwidth.

In this paper we advocate moving away from the common practice of implicitly assuming that all data exhibits spatial locality. We propose to reduce bandwidth consumption by fetching smaller blocks of data.

### 4. RUNAHEAD EXECUTION AND HARDWARE SCOUTING

Runahead execution has been proposed as a viable path to tolerating the ever increasing memory latencies. A runahead processor enters speculative runahead mode upon encountering a stall condition. Before runahead mode is entered, a checkpoint of the architected register state is taken. The execution then continues in the hope of finding independent memory operations further down the instruction stream. A memory operation is considered independent if its address computation is not dependent of the destination register of the load causing the launch into runahead mode or a previous cache miss. Throughout this paper we refer to the memory operation causing runahead execution as the *launch point*. Once the data of the launch point load miss has arrived, execution is resumed from the checkpointed register state. We will use the term re-execution to refer to normal mode execution of instructions that previously have been executed in runahead mode.

In runahead mode each register that is dependent on the destination register of the launch point load is marked by a Not-A-Thing bit. The effective address of memory operations whose source registers are not marked by the corresponding Not-A-Thing bit, can be computed and forwarded to the memory system as non binding prefetches. If such a prefetch misses in the cache, the destination register is marked by a Not-A-Thing bit. Similarly for stores where the effective address is known, a prefetch is forwarded to the memory system. Runahead can therefore be viewed as an intelligent prefetcher that operates when the processor would otherwise have been stalled.

Hardware Scouting, described by Chaudhry et. al.[7], is an extension of runahead execution, that includes several optimizations to previous runahead proposals. In hardware scouting, launching and exiting out of runahead is a zero-latency operation and runahead mode is also entered on low latency misses (L2 hits). This can be contrasted to a previous runahead proposal where runahead mode is entered first when an L2 miss has been detected[18]. In addition to entering runahead mode on a load miss, the hardware scouting proposal also identifies the case of a pending store when the store buffer is full as a stall event on which to launch into runahead mode.

To execute a dependent instruction, i.e. an instruction with a source register marked as Not-A-Thing, simply requires an *OR* operation of the Not-A-Thing bits associated with the source registers. Therefore, by providing specialized Not-A-Thing functional units, runahead mode enables a faster execution than normal mode since the instructions producing a Not-A-Thing value will not consume regular execution unit resources. In runahead mode the serialization enforced by certain instructions (e.g. CASA) can also be ignored, further speeding up the runahead mode execution.

## 5. INDEPENDENT SPATIAL CACHE MISSES

Throughout this paper we use the term *spatial miss*<sub>(S,C)</sub> to refer to a cache miss in a system with a fetch block size of  $S$  that would have been avoided if the larger cache block size  $C$  would have been used. To isolate from replacement algorithm effects we use the metric in the context of subblocked caches<sup>2</sup>[20], where  $C$  bytes of data map to the same cache tag and each subblock is of size  $S$ . An example of a *spatial miss*<sub>(16,64)</sub> is an access that leads to a cache miss in a cache with 64 byte cache blocks divided into 16 byte subblocks but that would hit if the cache was not subblocked.

The most common example of spatial misses occur when traversing an array, which often demonstrates excellent spatial locality. The address computation is often based on a base register holding the base address of the array and an offset based on a loop variable  $i$  that is not dependent on a memory access. Similarly for more complex data structures, compilers often use fixed offsets to access member variables. Since these accesses often are based on address computations that are not dependent on previous memory accesses, it leads us to the hypothesis explored in this paper:

*Spatial misses encountered in runahead mode are based on address computations of independent registers.*

Independent registers refers to source registers of a runahead mode memory operation that is not, directly or indirectly, dependent on any previous cache miss, including the launch point memory access. If we assume that a large fraction of the cache misses are encountered in runahead mode and that address computations of spatial misses are independent, most spatial misses would lead to prefetches in a runahead processor. If the assumption that many misses occur in runahead does not hold, it would imply that cache misses are so far apart that the runahead thread launched on a miss does not cover the next miss. If that is the case it is unlikely that the performance of such applications are limited by the memory system, and are therefore less interesting for this kind of study.

Sometimes in runahead mode the effective address of a memory operation can not be computed due to one of the source registers being marked as Not-A-Thing by a previous instruction. One example of this behavior is a linked-object list traversal. If the load of a list-item pointer misses in the cache, runahead execution will not be able to prefetch any of the subsequent list items since all items are dependent on the launch-point load. From a spatial access point of view, this case is not affected by the use of smaller cache blocks unless multiple objects fit in the same cache block. The reason for this is that the compiler usually uses register plus offset based addressing for this kind of accesses and if the register is Not-A-Thing, the accesses can not be sent out in either way. If we assume that all address computations of spatial misses are independent, all spatial misses will be found except for accesses of memory operations not covered by the runahead thread. One such example is if the runahead thread is terminated when only part of a cache block's spatial misses have been discovered.

<sup>2</sup>Sometimes also referred to as a sector cache

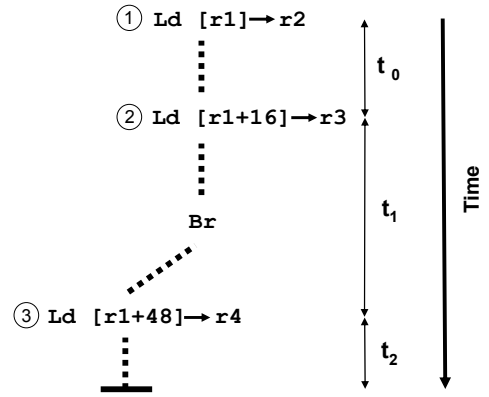


Figure 2: Temporal Aspects of Runahead Prefetching

It is worth noting that while runahead execution can be very effective at prefetching instructions into the instruction cache, runahead can not explore spatial instruction misses. The reason is that once an instruction miss occurs, the runahead thread can not continue and is therefore not able to find the spatial misses. Instruction misses are therefore costlier than data misses and presents a challenge since a runahead launch only can prefetch one cache block of instructions.

## 6. TEMPORAL ASPECTS

In a system where spatial misses lead to separate data fetches, there is a temporal aspect. Spatial misses occurring in runahead mode are requested later than they would be in a system with a larger cache block size. If we assume that re-execution in normal mode is equally fast as runahead execution, then due to the runahead way of returning to a checkpoint and re-executing instructions, the delayed fetch of spatial misses should not incur any extra delay. The reason for this is that by the time they are re-executed the same amount of time has passed as the latency of the launch point miss.

This is illustrated in Figure 2 where the first load into  $r2$  occurs in normal mode and triggers runahead execution. Then the second load into  $r3$  is a spatial miss<sub>(64,16)</sub> with respect to the first load, i.e. the launch point load. The second load is forwarded to the memory system  $t_0$  cycles after the launch point load access is sent out. However the re-execution will not occur until  $t_1 + t_2 + t_0$  cycles after it was fetched, which is equal to the latency of the launch point load access. The same reasoning holds if the first load is not a launch point miss, but instead a regular cache miss occurring in runahead mode.

As mentioned in Section 4 runahead execution can be faster than normal mode execution since additional specialized functional units can execute instructions marked as Not-A-Thing and serialization can be ignored. Therefore, despite the delayed request, the data have a high probability of being available in a timely fashion.

## 7. ARCHITECTURAL PROPOSAL

The implications of our observations from a memory system design point of view is that spatial misses in runahead to a much higher degree than before can be identified and sent out as prefetches. Therefore memory systems do not need to rely on the prefetching effect of large cache blocks and instead allow smaller pieces of data to be requested separately. To exploit this observation we propose a method called *fine-grained fetching*, which saves bandwidth by avoiding fetching of unnecessary data. Architecturally, this could be implemented by using a small cache block size in the highest level cache. To avoid the increase in tag overhead associated with small cache blocks, we evaluate our system using a subblocked cache organization.

The fine-grained fetch method will lead to an increased on-chip address bandwidth. While on-chip bandwidth also is a restricted resource we believe it to be a less constrained resource than off-chip bandwidth. The off-chip address bandwidth consumption will also increase and we will discuss in Section 10.1 how to attack this issue.

The runahead fine-grained fetch method is unsuitable to fetching instructions, since runahead can not find spatial instruction misses beyond an instruction miss. Since instructions often exhibit a high degree of spatial locality, we advocate the use of larger cache blocks for instruction caches. For unified second level caches, this leads to a complication due to the complexity involved in maintaining coherence between a L1 cache with a larger cache block size than the L2 cache. The solution proposed in this paper is to use a small L1 cache block size for both the L1 instruction and data cache, and automatically prefetch multiple continuous cache blocks for each instruction miss into the L1 instruction cache. More advanced discontinuous instruction prefetchers[23] have been proposed, but are not evaluated in this paper. It is worth noting that a successful use of such a predictor could result in a bandwidth reduction also for instructions, since the fetching of wrong-path instruction subblocks could be avoided.

## 8. METHODOLOGY

Our evaluation of runahead fine-grained fetch is based on full system simulation using both commercial and scientific workloads.

### 8.1 Simulation Setup

Simics is an execution-driven simulator that models a SPARC V9 system accurately enough to run Solaris unmodified[16]. To model the timing of a runahead processor we have extended Simics with a processor model, *Headrunner*, and a detailed memory hierarchy simulator VASA[9]. While equally relevant for single thread performance, we have opted to evaluate fine-grained fetching in chip multiprocessor context due to increased bandwidth requirement.

#### 8.1.1 Benchmarks

We use five multithreaded benchmarks. Two commercial applications SPECJBB 2000 and APACHE. SPECJBB 2000 (JBB2000) is a commercial JAVA-based middleware benchmark which evaluates the performance of server-side JAVA [24]. APACHE is a benchmark modeling the Apache open source

Web server to which URL-requests are sent by a client [1]. SPECJBB 2000 is run for 2000 transactions after a warm-up period of 100 000 transactions. APACHE was run for 800 transactions with a 1000 transactions warm-up period.

In addition we used three randomly chosen applications from the SPLASH2 benchmark suite, LU\_NC, WATER\_S and WATER\_N. We ran the SPLASH2 benchmarks with the default input sets specified in the SPLASH2 code and warmed the caches according to Woo et al [25].

#### 8.1.2 Processor and Memory System Configuration

We model a highly simplified non-pipelined, in-order and single-issue processor model. Hence, no wrong path effects are modeled. We believe that the exclusion of wrong-path effects in this study is likely to lead to an underestimation of the benefit of fine-grained fetch, since fetching smaller blocks of data along the wrong path is likely to have positive performance effects.

The CMP we are simulating contains 16 single-threaded cores each capable of invoking a runahead thread. Each core have separate L1 instruction and data caches, and a L2 cache which is shared among all the cores. We evaluate fine-grained fetch with a subblocked L2 cache configuration. To reduce interaction effects we have opted to simplify our memory system by assuming an infinite storebuffer and a perfect instruction cache.

Processor	16 cores per chip
Frequency	4 GHz
L1 data cache	64 KB, 8-way, 1 cycles
L1 block size	16 B
Shared L2 cache	2 MB, 8-way, 24 cycles
L2 block size	64 B subblocked 1, 2, 4 times
Memory Latency	200 cycles

Table 1: Simulated Target System Parameters

We model separate unidirectional address and data buses of varying bandwidth. We assume an address packet size of 8 bytes, write packet header of 4 bytes and data packet header of 4 bytes, where each bus packet includes control state and ECC. We do not model any additional bus latency instead only the queueing delay is added. Further details of the simulated configuration can be found in Table 1.

#### 8.1.3 The Headrunner Simulator

The Headrunner simulator is a simplified and idealized model of an in-order runahead processor. Runahead mode is entered on a data cache miss. We have assumed that a check-point can be made without incurring any stall time. The execution returns to normal mode when the cache block associated with the launch point load is filled. If the runahead thread comes across a memory access with side-effects, an asynchronous trap or an external interrupt, the execution stalls until relaunch into normal mode. Furthermore, if a miss is encountered in normal execution mode and the application is in kernel mode, the simulator does not enter runahead mode. The reason for this is simulator simplicity. Many simulation-wise complicated corner cases tend to happen when launching into runahead in kernel mode.

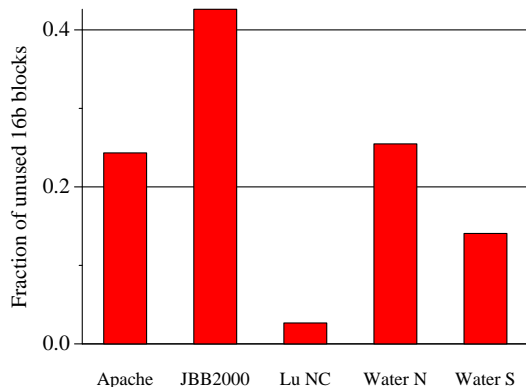


Figure 3: The fraction of unused 16 byte subblocks at eviction with a 64B cache block.

## 9. RESULTS

### 9.1 Spatial Usage

To measure the potential bandwidth savings achievable by using a smaller cache block size, we have measured the number of untouched 16 byte subblocks upon eviction, assuming a 64 byte cache block size. As can be seen in Figure 3 the untouched subblocks constitutes up to 40% of the allocated subblocks, showing a very low spatial utilization. Previous studies have reported that less than half of the allocated data gets used before eviction even for a small cache block size of 32 byte[15]. These pieces of data consume bandwidth without contributing to performance. In a bandwidth constrained system these extra subblocks lead to increased queuing delay and therefore increases the latency of other memory requests. If these blocks could be identified ahead of time, the amount of data transferred and therefore the memory access time, could be reduced.

We start our evaluating by using a subblocked L2 cache. We use subblocking instead of just reducing the cache block size to isolate the miss ratio effect from replacement effects. In a subblocked cache the same replacements will be made regardless of subblocking degree. Throughout this paper we refer to the number of subblocks mapping to the same tag as degree of subblocking, where a subblocking degree of one corresponds to a non-subblocked cache.

One of the potential obstacles towards successful runahead fine-grained fetching is if spatial misses can not be forwarded to the memory system due to source registers being dependent (marked as Not-A-Thing). Due to the way the head-runner simulator is implemented we can obtain the effective address also of memory operations where the source register is Not-A-Thing. We are therefore able to send non-intrusive probes to the memory hierarchy for Not-A-Thing memory operations and count the number of accesses that would have lead to a hit, block miss or a subblock miss. A subblock miss is an access that leads to a tag match but where the valid bit is not set for the subblock. Table 2 shows the number of Not-A-Thing misses in relation to (normal mode) instructions, L2 misses and L2 accesses.

Across all benchmarks the total number of Not-A-Thing misses are relatively few compared to the number of normal mode instructions. JBB2000 shows the highest amount

	NaT misses per 1000 instructions	NaT misses per L2 miss	NaT misses per L2 Access
Apache	2.04	0.07	0.01
JBB2000	13.52	0.73	0.14
Lu NC	0.00	0.00	0.00
Water N	0.06	0.13	0.00
Water S	0.00	0.00	0.00

Table 2: Number of Not-A-Thing misses.

which could be due to the object-oriented nature of the workload with pointer-indirections etc. JBB2000 also shows a very high number of Not-A-Thing misses per L2 miss fraction, indicating that the number of memory accesses that could not be forwarded to the memory system due to dependencies are high. From the perspective of fine-grained fetching, the question is however how many of these misses are spatial misses.

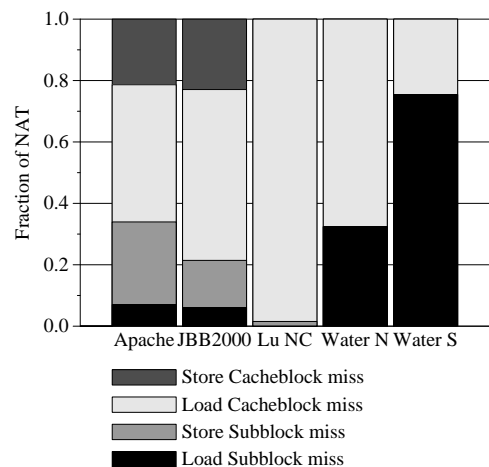
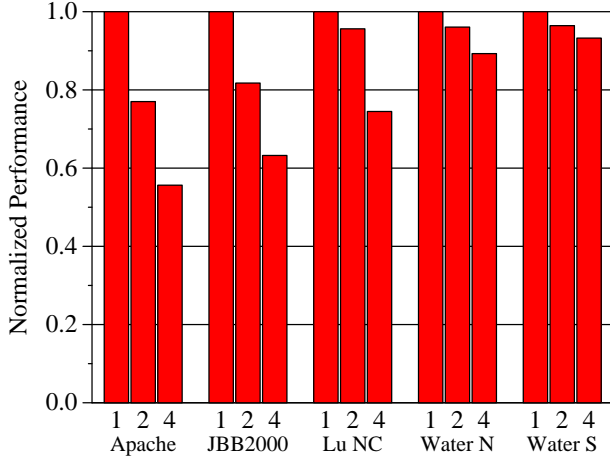


Figure 4: Distribution of memory accesses with Not-A-Thing marked source registers (64 B cache block 4 subblocks)

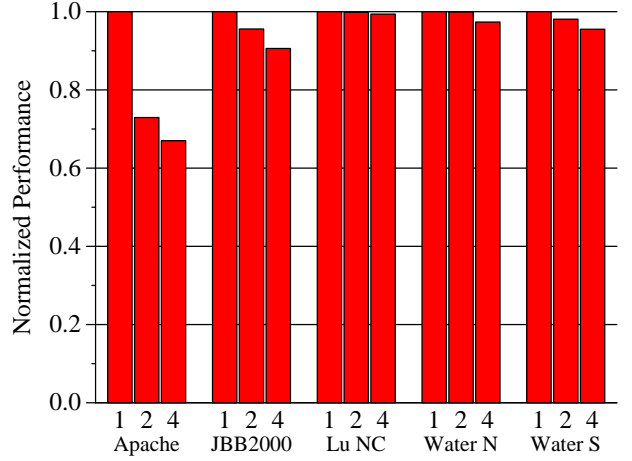
Figure 4 shows a breakdown of the Not-A-Thing misses into cache misses and subblock misses for both loads and stores. We find the fraction of load subblock misses to be low for both JBB2000 and APACHE. The store subblock misses are also fairly low although significantly higher than the load subblock miss fraction. This leads us to conclude that the high number of Not-A-Thing misses observed for JBB2000 in Table 2 includes relatively few spatial misses that would penalize the fine-grained fetch method. The breakdowns are less interesting for the SPLASH applications since they barely have any Not-A-Thing misses.

## 10. ISOLATING SUBBLOCK MISS PENALTY

To isolate the miss penalty effect we start off by ignoring bandwidth effects, i.e. a data fetch always takes a constant number of cycles regardless if 64 bytes or 16 bytes are fetched. Figure 5 a, shows the performance effect of increasing the degree of subblocking, i.e. using smaller and smaller subblocks, in a conventional in-order processor. The performance cost for fetching smaller subblocks is significant.



(a) In-order processor



(b) Runahead processor

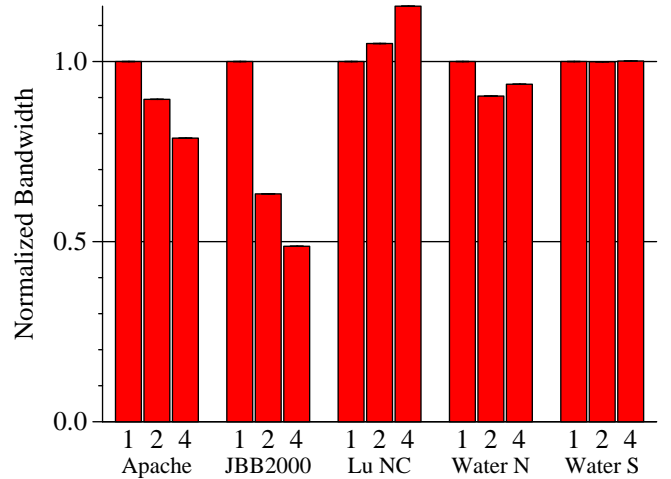
**Figure 5: Performance effect of reducing fetch block size in a system with infinite bandwidth. The 64B cache block is divided into 1, 2 and 4 subblocks**

Three out of the studied applications suffer an overall performance degradation of 25% or more. Especially the commercial workloads, APACHE and JBB2000 are negatively affected with up to 45% and 35% reduction respectively.

In an runahead processor, a cache miss can be overlapped by additional fetches in runahead mode. Due to the overlapping effect (memory level parallelism), the number of cache misses can not entirely indicate overall performance. As can be seen in Figure 5 b the isolated miss penalty is significantly lower than in the in-order case. Except for APACHE none of the applications suffer more than 10% performance degradation. In the case of APACHE the result is likely affected by the kernel mode limitation of the headrunner simulator. Since we have observed that 70% of L2 cache misses occur in kernel mode in APACHE, the simulator limitation of not launching into runahead mode when encountering a kernel mode stall condition, is likely to have a significant effect.

## 10.1 Taking bandwidth effects into account

Figure 5 b showed that from a miss ratio perspective the performance cost of decreasing fetch block size is limited for four out of the five studied applications. However reducing fetch block size also affects the bandwidth consumption. In Figure 6 we show the normalized data bandwidth when reducing the fetch block size. In the case of JBB2000 the bandwidth is reduced by 51%, while APACHE show a 21% reduction. The data bandwidth consumption of LU NC is increased by 15%, which can be explained by very high spatial locality in combination with additional data packet header overhead. The data packet header size of 4 byte which we have assumed leads to a four times as high overhead for 16 byte fetches compared to 64 byte fetches. This is somewhat exaggerated since the ECC does scale with the data size. We find that the bandwidth savings correlates well with Figure 3, since we see the largest effect in JBB2000 and the smallest effect for LU NC. Note that the bandwidth savings shown Figure 6 includes all fetched data, while the Figure 3 only



**Figure 6: Normalized data bandwidth when scaling the degree of subblocking for a cache with 64 B cache blocks**

measured spatial utilization on the evicted cache blocks.

Fine-grained fetching results in an address bandwidth increase in applications with a high spatial utilization. We have observed that the increase in address bandwidth, sometimes erases the savings made in data bandwidth. We have found that for the Splash applications, the high spatial utilization leaves few subblocks to be saved and can not offset the increase in address bandwidth. Therefore, in terms of total bandwidth, we only see a reduction for the commercial workloads, APACHE and JBB2000.

While beyond the scope of this paper, we believe multiple techniques can be applied to reduce the address overhead. We have observed that spatial misses often are encountered sequentially, indicating that bus encoding schemes may be

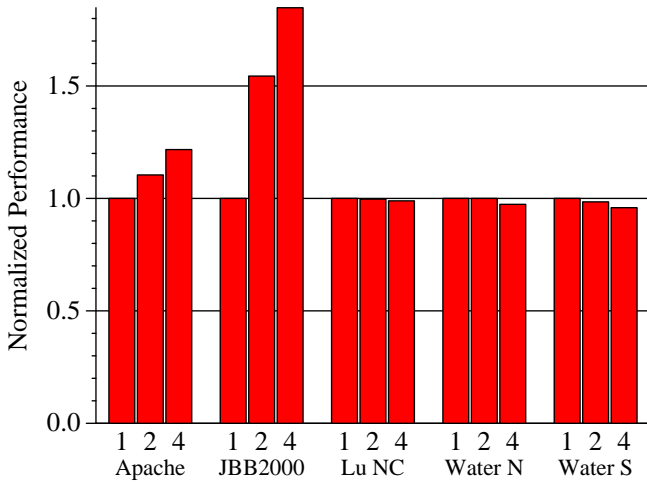


Figure 7: Performance Effect of increasing the degree of subblocking in a system with 8 GB/s data bandwidth

very successful at reducing the address overhead. Another possibility is to collapse multiple subblock requests into larger requests when requests are queued up on-chip.

By adding a bus model to our experimental setup we can evaluate the performance effect of conserving bandwidth through reduced fetch block size. Our model includes both an address and data bus, each providing the same bandwidth. Figure 7 shows the normalized performance in a system with 8 GB/s of dedicated data bandwidth (16 GB/s in total). This can be compared with the total (address and data) memory bandwidth of the Itanium 2 and Power4 processors, which are 6.4 GB/s and 12.8 GB/s respectively [19]. For the Power4 the bandwidth is divided equally in address and data bandwidth, hence the data bandwidth is 6.4 GB/s.

We can observe that for APACHE and JBB2000, the miss ratio penalty is outweighed by the bandwidth effect and reducing the fetch block size results in a speedup. For APACHE the speedup is 22% and for JBB2000 the speedup is 85%. The reason why we see such a large performance increase, despite a fairly modest bandwidth reduction, is likely due to burstiness in the memory requests. Since our bandwidth measurements are made for the entire run it does not reveal how much of the bandwidth is saved during the performance critical periods of contention.

For the Splash applications we note a very similar behavior as observed in Figure 5 b, indicating that the Splash applications are insensitive to bandwidth constraints. The reason is that the working sets of our benchmark setups are too small and fits in the relatively large L2 cache of 2 MB. We have nevertheless chosen to include them in this study in order to demonstrate the performance stability of the the fine-grained fetch method.

We have also simulated a system with 16 byte cache blocks, in order to estimate the performance increase possible from using small cache blocks instead of subblocking. We ignore the increase in tag overhead and simulate a 2MB L2

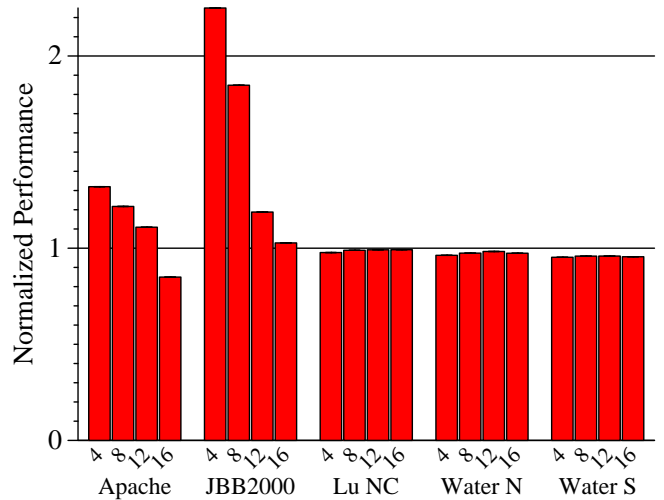


Figure 8: Normalized Performance when using a subblock size of 16 B (compared to 64 B) and scaling up the data bandwidth (4, 8, 12, 16 GB/s)

cache with 16 byte cache block and 8 GB/s of data bandwidth. For Apache and the Splash applications, the use of a smaller cache block size does not yield any significant performance improvement over the use of subblocking. However, for JBB2000 the speedup went from 85% with subblocking to 115% with small cache blocks.

In Figure 8 we show how performance is affected when scaling up the data bandwidth. As expected when we increase the bandwidth the positive performance effect decreases. But also if we decrease the bandwidth to 4 GB/s the speedup of APACHE and JBB2000 increases to 32% and 125% respectively. Since we are assuming a fixed DRAM access latency regardless of fetch size, the only performance contribution from smaller cache blocks are gained from a reduction in bus queuing delay. Therefore when we scale up the bandwidth the performance approaches the fixed latency without bandwidth modeling result from Figure 5 b. The interesting thing to note here is that we need to scale up the bandwidth fairly high in order for the performance improvement to diminish.

## 10.2 Power and Latency Implications of Fine-grained fetching

Chip power is one of the most constrained resources in processor design today. A significant source of power consumption are the inter-chip buses. Signal pin drivers can consume 15 to 20 percent of the total chip budget [17]. Therefore by reducing the amount of data transferred by using the fine-grained fetch method, substantial overall chip power savings can be made.

In narrow channel memory technologies like RAMBUS and DDR, a cache block request is divided into multiple reads (bursts) of e.g. 8 bytes as in DDR. If burst-scheduling is used, the bursts that make up a cache block, are read sequentially from a DRAM bank. In such systems the latency of reading a cache block is linear with the number of bursts. Therefore fine-grained fetching could also result in latency

reductions, further reducing the memory access cost.

## 11. RELATED WORK

In a patent by Burger and Wood[6] they propose an dynamic number of sub-block fetching where the number of sub-blocks fetched is determined by inspecting the usage of previously allocated blocks. They also propose fetching a pattern of discontinuous subblocks if an historical pattern is discovered. Kumar and Wilkerson proposes the spatial footprint predictor[15], that predicts how much data to fetch for each cache miss.

Several proposals have been made suggesting memory compression to save bandwidth[2][13][12]. We deem these approaches orthogonal to fine-grained fetch since compression of subblocks could further alleviate the bandwidth problem. Proximity communication[10] appears to be a very promising new technology. If proven successful it could very well disrupt current bandwidth trends.

Spracklen et. al. [22] presented latency and bandwidth implications of Hardware Scouting and describes several challenges involved in the design of chip-Multithreaded (CMT) processors. Other work targeted at understanding, exploiting or optimizing runahead execution have been presented by Sorin et. al.[14] and Chou et. al.[8].

## 12. CONCLUSIONS

Runahead execution is a very promising approach to tolerate long memory latencies. In this paper we have identified a characteristic of runahead execution that applies to spatial memory accesses. We have also demonstrated how this simple observation can be exploited to alleviate one of the most pressing memory system design challenges, the pin bandwidth of a chip. Our proposal of reducing bandwidth consumption by employing smaller fetch blocks have been shown to lead to execution time speedups of up to 125%.

This is the first step towards exploiting the characteristic of independent spatial misses through the fine-grained fetch method. However, we believe that future proposals can extend the method to reduce the increase in address bandwidth and reduce the performance penalty for workloads with a high degree of spatial utilization.

## 13. ACKNOWLEDGMENTS

I would like to thank Mark Hill and Kevin Moore for valuable comments on this paper and Håkan Zeffner for valuable suggestions on the Headrunner simulator implementation.

## 14. REFERENCES

- [1] A. R. Alameldeen, M.M.K. Martin, C.J Mauer, K.E. Moore, Xu Min, M.D. Hill, D.A Wood, and D.J. Sorin. Simulating a \$2m Commercial Server on a \$2K PC. *IEEE Computer*, 36(2):50–57, 2003.
- [2] Alaa R. Alameldeen and David A. Wood. Adaptive Cache Compression for High-Performance Processors. In *Proceedings of the 31th Annual International Symposium on Computer Architecture (ISCA'04)*, page 212, Washington, DC, USA, 2004. IEEE Computer Society.
- [3] Rajeev Balasubramonian, Sandhya Dwarkadas, and David H. Albonesi. Dynamically allocating processor resources between nearby and distant ilp. In *Proceedings of the 28th Annual International Symposium on Computer Architecture (ISCA'01)*, pages 26–37, 2001.
- [4] Ronald D. Barnes, Erik M. Nystrom, John W. Sias, Sanjay J. Patel, Nacho Navarro, and Wen mei W. Hwu. Beating in-order stalls with "flea-flicker" two-pass pipelining. In *MICRO 36: Proceedings of the 36th annual IEEE/ACM International Symposium on Microarchitecture*, page 387, 2003.
- [5] Doug Burger, James R. Goodman, and Alain Kgi. Limited bandwidth to affect processor design. *IEEE Micro*, 17(6), 1997.
- [6] Douglas C. Burger and David A. Wood. Cache with dynamic control of sub-block fetching. U.S. Patent No 6,557,080 issued April 29, 2003.
- [7] Shailender Chaudhry, Paul Caprioli, Sherman Yip, and Marc Tremblay. High-Performance Throughput Computing. *IEEE Micro*, 25(3):32–45, 2005.
- [8] Yuan Chou, Brian Fahs, and Santosh Abraham. Microarchitecture optimizations for exploiting memory-level parallelism. In *Proceedings of the 31th Annual International Symposium on Computer Architecture (ISCA'04)*, page 76, 2004.
- [9] Dan Wallin and Hkan Zeffner and Martin Karlsson and Erik Hagersten. VASA: A Simulator Infrastructure with Adjustable Fidelity. In *Proceedings of the International Conference on Parallel and Distributed Computing and Systems*, November 2005.
- [10] R.J. Drost, R.D. Hopkins, R. Ho, and I.E. Sutherland. Proximity communication. *IEEE Journal of Solid-State Circuits*, 39(9):1529– 1535, Sept 2004.
- [11] James Dundas and Trevor N. Mudge. Improving Data Cache Performance by Pre-Executing Instructions Under a Cache Miss. In *International Conference on Supercomputing*, pages 68–75, 1997.
- [12] Magnus Ekman and Per Stenström. A robust main-memory compression scheme. In *Proceedings of the 32nd Annual International Symposium on Computer Architecture (ISCA'05)*, pages 74–85, 2005.
- [13] Erik G. Hallnor and Steven K. Reinhardt. A unified compressed memory hierarchy. In *HPCA '05: Proceedings of the 11th International Symposium on High-Performance Computer Architecture*, pages 201–212, 2005.
- [14] Sorin Iacobovici, Lawrence Spracklen, Sudarshan Kadambi, Yuan Chou, and Santosh G. Abraham. Effective stream-based and execution-based data prefetching. In *ICS '04: Proceedings of the 18th annual international conference on Supercomputing*, pages 1–11, 2004.



- [15] Sanjeev Kumar and Christopher Wilkerson. Exploiting spatial locality in data caches using spatial footprints. In *Proceedings of the 25th Annual International Symposium on Computer Architecture (ISCA'98)*, pages 357–368, 1998.
- [16] "P. S. Magnusson, M. Christensson, D. Forsgren, J. Eskilson, G. Hillberg, J. Hgberg, A. Moestedt, F. Larsson, and B. Werner". Simics: A Full System Simulation Platform. *IEEE Computer*, pages 50–58, February 2002.
- [17] Trevor Mudge. Power: A First-Class Architectural Design Constraint. *IEEE Computer*, 34(4), April 2001.
- [18] O. Mutlu, J. Stark, C. Wilkerson, and Y. Patt. Runahead execution: An alternative to very large instruction windows for out-of-order processors. In *Proceedings of the Ninth International Symposium on High Performance Computer Architecture (HPCA-9)*, 2003.
- [19] Microprocessor Report Newsletter. Power5 Tops On Bandwidth, December 2003.
- [20] Jeffrey B. Rothman and Alan Jay Smith. The pool of subsectors cache design. In *ICS '99: Proceedings of the 13th international conference on Supercomputing*, pages 31–42, 1999.
- [21] Semiconductor Industry Association (SIA). International Technology Roadmap for Semiconductors 2004 Update.
- [22] Lawrence Spracklen and Santosh G. Abraham. Chip multithreading: Opportunities and challenges. In *HPCA '05: Proceedings of the 11th International Symposium on High-Performance Computer Architecture*, 2005.
- [23] Lawrence Spracklen, Yuan Chou, and Santosh G. Abraham. Effective instruction prefetching in chip multiprocessors for modern commercial applications. In *HPCA '05: Proceedings of the 11th International Symposium on High-Performance Computer Architecture*, 2005.
- [24] <http://www.spec.org/osg/jbb2000/>.
- [25] Steven Cameron Woo, Moriyoshi Ohara, Evan Torrie, Jaswinder Pal Singh, and Anoop Gupta. The SPLASH-2 Programs: Characterization and Methodological Considerations. In *Proceedings of the 22nd Annual International Symposium on Computer Architecture (ISCA'95)*, 1995.