

# Run-time Spatial Locality Detection and Optimization

Teresa L. Johnson   Matthew C. Merten   Wen-mei W. Hwu  
Center for Reliable and High-Performance Computing  
University of Illinois, Urbana-Champaign, IL 61801  
{tjohnson,merten,hwu}@crhc.uiuc.edu

## Abstract

*As the disparity between processor and main memory performance grows, the number of execution cycles spent waiting for memory accesses to complete also increases. As a result, latency hiding techniques are critical for improved application performance on future processors. We present a microarchitecture scheme which detects and adapts to varying spatial locality, dynamically adjusting the amount of data fetched on a cache miss. The Spatial Locality Detection Table, introduced in this paper, facilitates the detection of spatial locality across adjacent cached blocks. Results from detailed simulations of several integer programs show significant speedups. The improvements are due to the reduction of conflict and capacity misses by utilizing small blocks and small fetch sizes when spatial locality is absent, and the prefetching effect of large fetch sizes when spatial locality exists.*

## 1 Introduction

This paper introduces an approach to solving the growing memory latency problem [1] by intelligently exploiting *spatial locality*. Spatial locality refers to the tendency for neighboring memory locations to be referenced close together in time. Traditionally there have been two main approaches used to exploit spatial locality. The first approach is to use larger cache blocks, which have a natural prefetching effect. However, large cache blocks can result in wasted bus bandwidth and poor cache utilization, due to fragmentation and underutilized cache blocks. Both negative effects occur when data with little spatial locality is cached. The second common approach is to prefetch multiple blocks into the cache. However, prefetching is only beneficial when the prefetched data is accessed in cache, otherwise the prefetched data may displace more useful data from the cache, in addition to wasting bus bandwidth. Particularly when using large block sizes, the amount of prefetching is fixed. However, the spa-

tial locality, and hence the optimal prefetch amount, varies across and often within programs.

As the available chip area increases, it is meaningful to spend more resources to allow intelligent control over latency-hiding techniques, adapting to the variations in spatial locality. For numeric programs there are several known compiler techniques for optimizing data cache performance. In contrast, integer (non-numeric) programs often have irregular access patterns that the compiler cannot detect and optimize. For example, the temporal and spatial locality of linked list elements and hash table data are often difficult to determine at compile time. This paper focuses on cache performance optimization for integer programs. While we focus our attention on data caches, the techniques presented here are applicable to instruction caches.

In order to increase data cache effectiveness for integer programs we are investigating methods of *adaptive cache hierarchy management*, where we intelligently control caching decisions based on the usage characteristics of accessed data. In this paper we examine the problem of detecting spatial locality in accessed data, and automatically control the fetch of multiple smaller cache blocks into all data caches and buffers. Not only are we able to reduce the conflict and capacity misses with smaller cache lines and fetch sizes when spatial locality is absent, but we also reduce cold start misses and prefetch useful data with larger fetch sizes when spatial locality is present.

We introduce a new hardware mechanism called the *Spatial Locality Detection Table (SLDT)*. Each SLDT entry tracks the accesses to multiple adjacent cache blocks, facilitating detection of spatial locality across those blocks while they are cached. The resulting information is later recorded in the *Memory Address Table* [2] for long-term tracking of larger regions called *macroblocks*. We show that these extensions to the cache microarchitecture significantly improve the performance of integer applications, achieving up to 17% and 26% improvements for 100 and 200-cycle memory latencies, respectively. This scheme is fully compatible with existing Instruction Set Architectures (ISA).

The remainder of this paper is organized as follows: Section 2 discusses related work; Section 3 discusses general spatial locality issues; Section 4 discusses hardware techniques; Section 5 presents simulation results; and Section 6 concludes with future directions.

---

Copyright ©1997 IEEE. Published in the Proceedings of Micro-30, December 1-3, 1997 in Research Triangle Park, North Carolina. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works, must be obtained from the IEEE. Contact: Manager, Copyrights and Permissions / IEEE Service Center / 445 Hoes Lane / P.O. Box 1331 / Piscataway, NJ 08855-1331, USA. Telephone: + Intl. 908-562-3966.

## 2 Related Work

Several studies have examined the performance effects of cache block sizes [3][4]. One of the studies allowed multiple consecutive blocks to be fetched with one request [3], and found that for data caches the optimal statically-determined fetch size was generally twice the block size. In this work we also examine fetch sizes larger than the block size, however, we allow the fetch size to vary based on the detected spatial locality. Another method allows the number of blocks fetched on a miss to vary across program execution, but not across different data [5].

Hardware [6][7][8][9][10] and software [11][12][13] prefetching methods for uniprocessor machines have been proposed. However, many of these methods focus on prefetching regular array accesses within well-structured loops, which are access patterns primarily found in numeric codes. Other methods geared towards integer codes [14][15] focus on compiler-inserted prefetching of pointer targets, and could be used in conjunction with our techniques.

The dual data cache [16] attempts to intelligently exploit both spatial and temporal locality, however the temporal and spatial data must be placed in separate structures, and therefore the relative amounts of each type of data must be determined a priori. Also, the spatial locality detection method was tuned to numeric codes with constant stride vectors. In integer codes, the spatial locality patterns may not be as regular. The split temporal/spatial cache [17] is similar in structure to the dual data cache, however, the runtime locality detection mechanism is quite different than that of both the dual data cache and this paper.

## 3 Spatial Locality

Caches seek to exploit the principle of locality. By storing a referenced item, caches exploit *temporal locality* - the tendency for that item to be referenced soon. Additionally, by storing multiple items adjacent to the referenced item, they exploit *spatial locality* - the tendency for neighboring items to be referenced soon. While exploitation of temporal locality can result in cache hits for future accesses to a particular item, exploitation of spatial locality can result in cache hits for future accesses to multiple nearby items. Traditionally, exploitation of spatial locality is achieved through either larger block sizes or prefetching of additional blocks.

For a 32-byte cache block, we found that over half the time the extra data fetched into the cache simply wasted bus bandwidth and cache space [18]. Therefore, it would be beneficial to tune the amount of data fetched and cached on a miss to the spatial locality available in the data. Also, we found routines in benchmarks such as SPEC92 *gcc* where the amount of spatial locality in data fetched by a particular load instruction varied depending on the function arguments [18]. As such, neither static analysis (if even possible) nor profiling will result in definitive or accurate spatial locality information for the load instructions. Dynamic analysis of the spatial locality in the data offers greater promise. Also, dynamic schemes do not require profiling, which many users are unwilling to perform, or ISA changes.

## 4 Techniques

### 4.1 Overview of Prior Work

In this section we briefly overview the concept of a *macroblock*, as well as the *Memory Address Table (MAT)*, introduced in an earlier paper [2] and utilized in this work.

We showed that cache bypassing decisions could be effectively made at run-time, based on the previous usage of the memory address being accessed. Other bypassing schemes include [19][20][16][21]. In particular, our scheme dynamically kept track of the accessing frequencies of memory regions called macroblocks. The macroblocks are statically-defined blocks of memory with uniform size, larger than the cache block size. The macroblock size should be large enough so that the total number of accessed macroblocks is not excessively large, but small enough so that the access patterns of the cache blocks contained within each macroblock are relatively uniform. It was determined that 1K-byte macroblocks provide a good cost-performance tradeoff.

In order to keep track of the macroblocks at run time we use an MAT, which ideally contains an entry for each macroblock, and is accessed with a macroblock address. To support dynamic bypassing decisions, each entry in the table contains a saturating counter, where the counter value represents the frequency of accesses to the corresponding macroblock. For details on the MAT bypassing scheme see [2]. Also introduced in that paper was an optimization geared towards improving the efficiency of L1 bypasses, by tracking the spatial locality of bypassed data using the MAT, and using that information to determine how much data to fetch on an L1 bypass. In this paper we introduce a more robust spatial locality detection and optimization scheme using the SLDT, which enables much more efficient detection of spatial locality. Our new scheme also supports fetching varying amounts of data into both levels of the data cache, both with and without bypassing. In practice this spatial locality optimization should be performed in combination with bypassing, in order to achieve the best possible performance, as well as to amortize the cost of the MAT hardware. The cost of the combined hardware is addressed elsewhere [18] due to space constraints.

### 4.2 Support for Varying Fetch Sizes

The varying fetch size optimization could be supported using subblocks. In that case the block size is the largest fetch size and the subblock size is  $\text{gcd}(\text{fetch\_size}_0, \dots, \text{fetch\_size}_n)$ , where  $n$  is the number of fetch sizes supported. Currently, we only support two power-of-two fetch sizes for each level of cache, so the subblock size is simply the smaller fetch size. However, the cache lines will be underutilized when only the smaller size is fetched.

Instead, we use a cache with small lines, equal to the smaller fetch size, and optionally fill in multiple, consecutive blocks when the larger fetch size is chosen. This approach is similar to that used in some prefetching strategies [22]. As a result, the cache can be fully utilized, even when the smaller sizes are fetched. It also eliminates conflict misses resulting from accesses to different subblocks. However, this approach

makes detection of spatial reuses much more difficult, as will be described in Section 4.3. Also, smaller block sizes increase the tag array cost [18]. In our scheme, the  $max\_fetch\_size$  data is always aligned to  $max\_fetch\_size$  boundaries. As a result, our techniques will fetch data on either side of the accessed element, depending on the location of the element within the  $max\_fetch\_size$  block. In our experience, spatial locality in the data cache can be in either direction (spatially) from the referenced element.

### 4.3 Spatial Locality Detection Table

To facilitate spatial locality tracking, a *spatial counter*, or *sctr*, is included in each MAT entry. The role of the *sctr* is to track the medium to long-term spatial locality of the corresponding macroblock, and to make fetch size decisions, as will be explained in Section 4.4. This counter will be incremented whenever a *spatial miss* is detected, which occurs when portions of the same larger fetch size block of data reside in the cache, but not the element currently being accessed. Therefore, a hit might have occurred if the larger fetch size was fetched, rather than the smaller fetch size. In our implementation, where multiple cache blocks are filled when the larger fetch size is chosen, a spatial miss is not trivial to detect. If the cache is not fully-associative, the tags for different blocks residing in the same larger fetch size block will lie in consecutive sets. Searching for other cache blocks in the same larger fetch size block of data will require access to the tags in these consecutive sets, and thus either additional cycles to access, or additional hardware support. One possibility is a restructured tag array design allowing efficient access to multiple consecutive sets of tags. Alternatively, a separate structure can be used to detect this information, which is the approach investigated in this work.

This structure is called the *Spatial Locality Detection Table (SLDT)*, and is designed for efficient detection of spatial reuses with low hardware overhead. The role of the SLDT is to detect spatial locality of data while it is in the cache, for recording in the MAT when the data is displaced. The SLDT is basically a tag array for blocks of the larger fetch size, allowing single-cycle access to the necessary information. Figure 1 shows an overview of how the SLDT interacts with the MAT and L1 data cache, where the double-arrow line shows the correspondence of four L1 data cache entries with a single SLDT entry. In order to track all cache blocks, the SLDT would need  $N$  entries, where  $N$  is the number of blocks in the cache. This represents the worst case of having fetched only smaller (line) size blocks into the cache, all from different larger size blocks. However, in order to reduce the hardware overhead of the SLDT, we use a much smaller number of entries, which will allow us to capture only the shorter-term spatial reuses. The same SLDT could be used to track the spatial locality aspects of all structures at the same level in the memory hierarchy, such as the data cache, the instruction cache, and, when we perform bypassing, the bypass buffer.

The SLDT tags correspond to maximum fetch size blocks. The *sz* field is one bit indicating if either the larger size block was fetched into the cache, or if only smaller blocks were fetched. The *vc* (*valid count*) field is

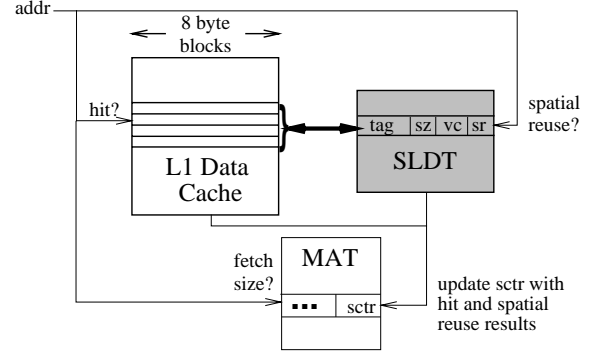


Figure 1: SLDT and MAT Hardware

$\log(max\_fetch\_size/min\_fetch\_size)$  bits in length, and indicates how many of the smaller blocks in the larger size block are currently valid in the data cache. The actual number of valid smaller blocks is  $vc+1$ . An SLDT entry will only be valid for a larger size block when some of its constituent blocks are currently valid in the data cache. A bit mask could be used to implement the *vc*, rather than the counter design, to reduce the operational complexity. However, for large maximum to minimum fetch size ratios, a bit mask will result in larger entries. Finally, the *sr* (*spatial reuse*) bit will be set if spatial reuse is detected, as will be discussed later.

When a larger size block of data is fetched into the cache, an SLDT entry is allocated (possibly causing the replacement of an existing entry) and the values of *sz* and *vc* are set to 1 and  $max\_fetch\_size/min\_fetch\_size - 1$ , respectively. If a smaller size block is fetched and no SLDT entry currently exists for the corresponding larger size block, then an entry is allocated and *sz* and *vc* are both initialized to 0. If an entry already exists, *vc* is incremented to indicate that there is now an additional valid constituent block in the data cache. For both fetch sizes the *sr* bit is initialized to 0. When a cache block is replaced from the data cache, the corresponding SLDT entry is accessed and its *vc* value is decremented if it is greater than 0. If *vc* is already 0, then this was the only valid block, so the SLDT entry is invalidated. When an SLDT entry is invalidated its *sr* bit is checked to see if there was any spatial reuse while the data was cached. If not, the corresponding entry in the MAT is accessed and its *sctr* is decremented, effectively depositing the information in the MAT for longer-term tracking. Because the SLDT is managed as a cache, entries can be replaced, in which case the same actions are taken.

An *fi* (*fetch initiator*) bit is added to each data cache tag to help detect spatial hits. The *fi* bit is set to 1 during the cache refill for the cache block containing the referenced element (i.e. the cache block causing the fetch), otherwise it is reset to 0. Therefore, a hit to any block with a 0 *fi* bit is a spatial hit, as this data was fetched into the cache by a miss to some other element.

Table 1 summarizes the actions taken by the SLDT for memory accesses. The *sr* bit, which was initialized to zero, is set for all types of both spatial misses and spatial hits. Two types of spatial misses are detected. The first type of spatial miss occurs when other portions of the same larger

Cache Access	SLDT Access	<i>fi</i>	<i>sz</i>	<i>vc</i>	Action
miss	hit	-	0		$sr = 1; sctr++$
		-	1		$sr = 1$
hit	hit	0			$sr = 1$
		0	>0		$sr = 1$
hit	miss	0	-	-	alloc SLDT entry; $sz = 1; sr = 1$
		1	-	-	alloc SLDT entry
Cache entry replaced		$vc > 0$		$vc--$	
		$vc == 0$		invalidate SLDT entry	
SLDT entry replaced or invalidated		$sr == 0$		$sctr--$	
		$sr == 1$		no action	

Table 1: SLDT Actions. A dash indicates that there is no corresponding value, and a blank indicates that the value does not matter.

fetch size block were fetched independently, indicated by a valid SLDT entry with a *sz* of 0. Therefore, there might have been a cache hit if the larger size block was fetched, so the corresponding entry in the MAT is accessed and its *sctr* is incremented. The second type can occur when the larger size block was fetched, but one of its constituent blocks was displaced from the cache, as indicated by a cache miss and a valid SLDT entry with a *sz* of 1. It is not trivial to detect if this miss is to the element which caused the original fetch, or to some other element in the larger fetch size block. The *sr* bit is conservatively set, but the *sctr* in the corresponding MAT entry is not incremented.

A spatial hit can occur in two situations. If the larger size block was fetched, then the *fi* bit will only be set for one of the loaded cache blocks. A hit to any of the loaded cache blocks without the *fi* bit set is a spatial hit, as described earlier. We do not increment the *sctr* on spatial hits, because our fetch size was correct. We only update the *sctr* when the fetch size should be changed in the future. When multiple smaller blocks were fetched, a hit to one of these is also characterized as a spatial hit. This case is detected by checking if *vc* is larger than 0 when *sz* is 0. However, we do not increment the *sctr* in this case either because a spatial miss would have been detected earlier when a second element in the larger fetch size block was first accessed (and missed).

#### 4.4 Fetch Size Decisions

On a memory access, a lookup in the MAT of the corresponding macroblock entry is performed in parallel with the data cache access. If an entry is found, the *sctr* value is compared to some threshold value. The larger size is fetched if the *sctr* is larger than the threshold, otherwise the smaller size is fetched. If no entry is found, a new entry is allocated and the *sctr* value is initialized to the threshold value, and the larger fetch size is chosen. In this paper the threshold is 50% of the maximum *sctr* value.

## 5 Experimental Evaluation

### 5.1 Experimental Environment

We simulate ten benchmarks, including *026.compress*, *072.sc* and *085.cc1* from the *SPEC92* benchmark suite using the *SPEC* reference inputs, and *099.go*, *147.vortex*, *130.li*, *134.perl*, and *124.m88ksim* from the *SPEC95* benchmark

L1 Icache	32K-byte split-block, direct mapped, 64-byte block
L1 Dcache	16K-byte non-blocking (50 max), direct mapped, 32-byte block, multiported, writeback, no write alloc
L1-L2 Bus	8-byte bandwidth, split-transaction, 4-cycle latency, returns critical word first
L2 Dcache	same as L1 Dcache except: 256K-byte, 64-byte block
System Bus	same as L1-L2 Bus except: 100-cycle latency
Issue	8-issue uniform, except 4 memory ops/cycle max
Registers	64 integer, 64 double precision floating-point

Table 2: Base Configuration.

suite using the training inputs. The last two benchmarks consist of modules from the *IMPACT* compiler [23] that we felt were representative of many real-world integer applications. *Pcode*, the front end of *IMPACT*, is run performing dependence analysis with the internal representation of the *combine.c* file from GNU CC as input. *lmdes2\_customizer*, a machine description optimizer, is run optimizing the SuperSPARC machine description. These optimizations operate over linked list and complex data structures, and utilize hash tables for efficient access to the information.

In order to provide a realistic evaluation of our technique for future high-performance, high-issue rate systems, we first optimized the code using the *IMPACT* compiler [23]. Classical optimizations were applied, then optimizations were performed which increase instruction level parallelism. The code was scheduled, register allocated and optimized for an eight-issue, scoreboarded, superscalar processor with register renaming. The ISA is an extension of the HP PA-RISC instruction set to support compile-time speculation.

We perform cycle-by-cycle emulation-driven simulation on a Hewlett-Packard *PA-RISC 7100* workstation, modelling the processor and the memory hierarchy (including all related busses). The instruction latencies used are those of a Hewlett-Packard *PA-RISC 7100*. The base machine configuration is described in Table 2.

Since simulating the entire applications at this level of detail would be impractical, uniform sampling is used to reduce simulation time [24], however emulation is still performed between samples. The simulated samples are 200,000 instructions in length and are spaced evenly every 20,000,000 instructions, yielding a 1% sampling ratio. For smaller applications, the time between samples is reduced to maintain at least 50 samples (10,000,000 instructions). To evaluate the accuracy of this technique, we simulated several configurations both with and without sampling, and found that the improvements reported in this paper are very close to those obtained by simulating the entire application.

### 5.2 Macroblock Spatial Locality Variations

Before presenting the performance improvements achieved by our optimizations, we first examine the accuracy of the macroblock granularity for tracking spatial locality. It is important to have accurate spatial locality information in the MAT for our scheme to be successful. This means that all data elements in a macroblock should have similar amounts of spatial locality at each phase of program execution.

After dividing main memory into macroblocks, as described in Section 4.1, the macroblocks can be further subdivided into smaller sections, each the size of a 32-byte cache block. We will simply call these smaller sections *blocks*. In

order to determine the dynamic cache block spatial locality behavior, we examined the accesses to each of these blocks, gathering information twice per simulation sample, or every 100,000 instructions. At the end of each 100,000-instruction phase, we determined the fraction of times that each block in memory had at least one spatial reuse each time it was cached during that phase. We call this the *spatial reuse fraction* for that block. Figure 2 shows a graphical representation of the resulting information for two programs. Each row in the graph represents a 1K-byte macroblock accessed in a particular phase. For every phase in which a particular macroblock was accessed, there will be a corresponding row. Each row contains one data point for every 32-byte block accessed during the corresponding phase that lies in that macroblock. For the purposes of clarity, the rows were sorted by the average of the block spatial reuse fractions per macroblock. The averages increase from the bottom to the top of the graphs. The cache blocks in each macroblock were also sorted so that their spatial reuse fractions increase from left to right. Some rows are not full, meaning that not all of their blocks were accessed during the corresponding phase. Finally, the cache blocks with spatial reuse fractions falling within the same range were plotted with the same marker.

Figure 2(a) shows the spatial locality distribution for *026.compress*. Most of the blocks, corresponding to the lighter gray points, have spatial reuse fractions between 0 and 0.25, meaning that there was spatial reuse to those blocks less than 25% of the time they were cached. Very few of the blocks, corresponding to the black points, had spatial reuse more than 75% of the time they were cached. This represents a fairly optimal scenario, because most of the macroblocks contain blocks which have approximately the same amount of reuse. Figure 2(b) shows the distribution for *134.perl*. Around 34% of the macroblocks (IDs 0 to 6500) contain only blocks with little spatial reuse, their spatial reuse fractions all less than 0.25. About 29% of the macroblocks (IDs 13500 to 18900) contain only blocks with large fractions of spatial reuse, their spatial reuse fractions all over 0.75. About 37% of the macroblocks contain cache blocks with differing amounts of spatial reuse. The medium gray points in some of these rows correspond to blocks with spatial reuse fractions between 0.25 and 0.75. However, this information does not reveal the time intervals over which the spatial reuse in these blocks varies. It is possible that in certain small phases of program execution the spatial locality behavior is uniform, but that it changes drastically from one small phase of execution to another. This type of behavior is possible due to dynamically-allocated data, where a particular section of memory may be allocated as one type of data in one part of the program, then freed and reallocated as another type later.

### 5.3 Performance Improvements

In this section we examine the performance improvement, or the execution cycles eliminated, over the base 8-issue configuration described in Section 5.1. To support varying fetch sizes, we use an SLDT and an MAT at each level of the cache hierarchy. The L1 and L2 SLDTs are direct-mapped with 32 entries. A large number of simulations showed that

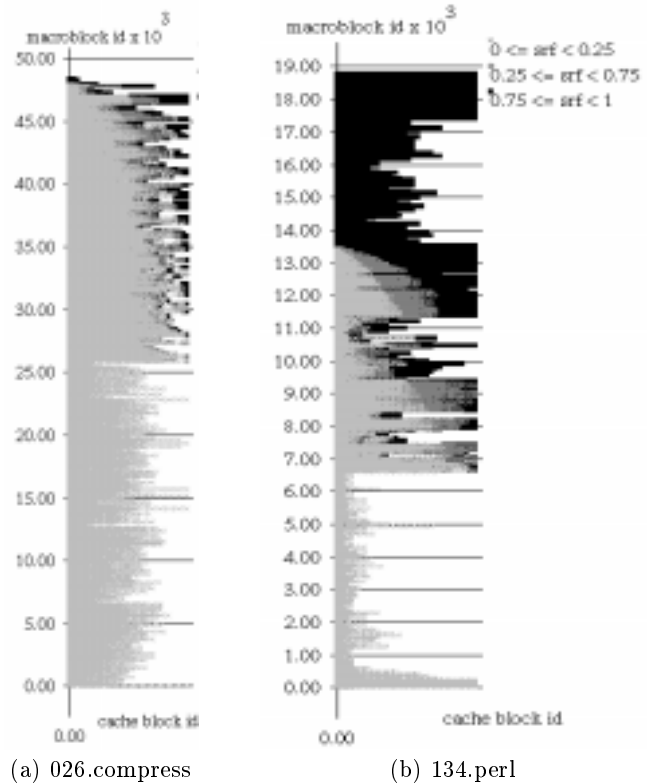


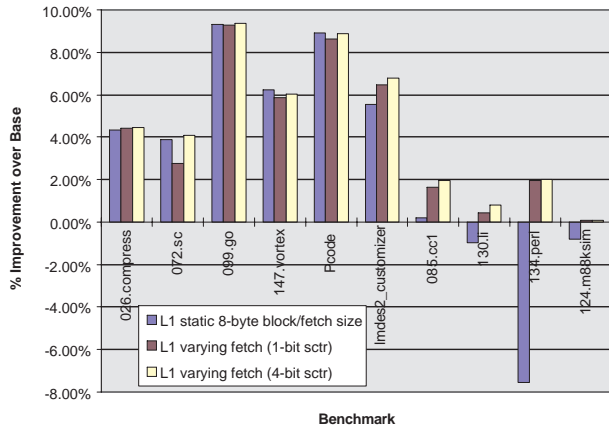
Figure 2: Spatial reuse fractions (*srf*) for cache-block-sized-data in the accessed macroblocks for two applications.

direct-mapped SLDTs perform as well as a fully-associative design, and that 32 entries perform almost as well as any larger power-of-two number of entries up to 1024 entries, which was the maximum size examined. The L1 and L2 MATs utilize 1K-byte macroblocks, and we examine both one and four-bit *sctrs*. We first present results for infinite-entry MATs, then study the effects of limiting the number of MAT entries.

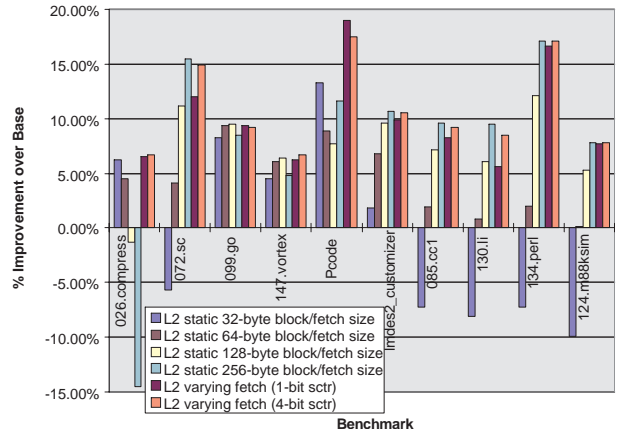
#### 5.3.1 Static versus Varying Fetch Sizes

The left bar for each benchmark in Figure 3(a) shows the performance improvement achieved by using 8-byte L1 data cache blocks with a static 8-byte fetch size, over the base 32-byte block and fetch sizes. These bars show that the better choice of block size is highly application-dependent. The right bars show the improvement achieved by our spatial locality optimization at the L1 level only, using an 8-byte L1 data cache block size, and fetching either 8 or 32-bytes on an L1 data cache miss, depending on the value of the corresponding *sctr*. The results show that our scheme is able to obtain either almost all of the performance, or is able to outperform, the best static fetch size scheme. In most cases the 1 and 4-bit *sctrs* perform similarly, but in one case the 4-bit *sctr* achieves almost 2% greater performance improvement.

The four leftmost bars for each benchmark in Figure 3(b) show the performance improvement using different L2 data cache block and (static) fetch sizes, and our L1 spatial lo-



(a) L1 Trends



(b) L2 Trends (with L1 varying fetches)

Figure 3: Performance for various statically-determined block/fetch sizes and for our spatial locality optimizations using both 1 and 4-bit *sctrs*.

cality optimization with a 4-bit *sctr*. The base configuration is again the configuration described in Section 5.1, which has 64-byte L2 data cache block and fetch sizes. These bars show that, again, the better static block/fetch size is highly application-dependent. For example, *134.perf* achieves much better performance with a 256-byte fetch size, while *026.compress* achieves its best performance with a 32-byte fetch size, obtaining over 14% performance degradation with 256-byte fetches. The rightmost two bars in Figure 3(b) show the performance improvement achieved with our L2 spatial locality optimization, which uses a 32-byte L2 data cache block size and fetches either 32 or 256 bytes on an L2 data cache miss, depending on the value of the corresponding L2 MAT *sctr*. Again, our spatial locality optimizations are able to obtain almost the same or better performance than the best static fetch size scheme for all benchmarks.

Figure 4 shows the breakdown of processor stall cycles attributed to different types of data cache misses, as a percentage of the total base configuration execution cycles. The left and right bars for each benchmark are the stall cycle breakdown for the base configuration and our spatial locality optimization, respectively. The spatial locality optimizations were performed at both cache levels, using the same configuration as in Figure 3(b) with a 4-bit *sctr*. For the benchmarks that have large amounts of spatial locality, as indicated from the results of Figure 3, we obtain large reductions in L2 cold start stall cycles by fetching 256 bytes on L2 cache misses. The benchmarks with little spatial locality in the L1 data cache, such as *026.compress* and *Pcode*, obtained reductions in L1 capacity miss stall cycles from fetching fewer small cache blocks on L1 misses. In some cases the L1 cold start stall cycles increase, indicating that the L1 optimizations are less aggressive in terms of fetching more data, however these increases are generally more than compensated by reductions in other types of L1 stall cycles. The L1 conflict miss stall cycles increase for *lmdes2\_customizer*, because it tends to fetch fewer blocks on an L1 miss, exposing some conflicts that were interpreted as capacity misses in the base configuration.

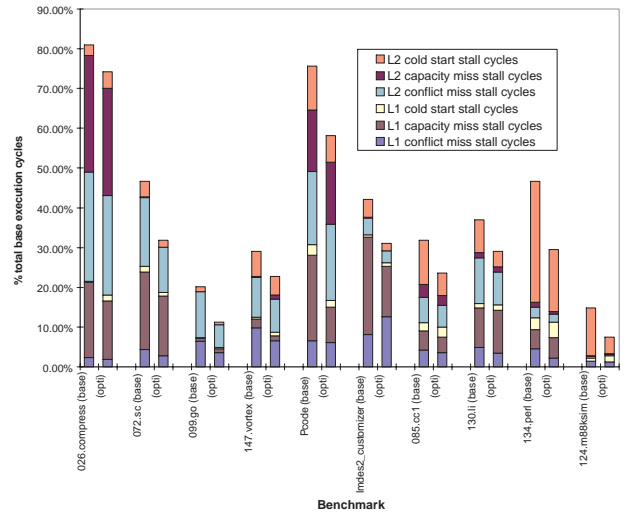


Figure 4: Stall cycle breakdown for base and the spatial locality optimizations.

### 5.3.2 Set-associative Data Caches

Increasing the set-associativity of the data caches can reduce the number of conflict misses, which may in turn reduce the advantage offered by our optimizations. However, the reductions in capacity and cold start stall cycles that our optimizations achieve should remain. To investigate these effects, the data cache configuration discussed in Section 5.1 was modified to have a 2-way set-associative L1 data cache and a 4-way set-associative L2 data cache.

Figure 5 shows the new performance improvements for our optimizations. The left bars show the result of applying our optimizations to the L1 data cache only, and the right bars show the result of applying our techniques to both the L1 and L2 data caches, using four-bit *sctrs*. The improvements have reduced significantly for some benchmarks over those shown in Figure 3. However, large improvements are still achieved for some benchmarks, particularly when applying

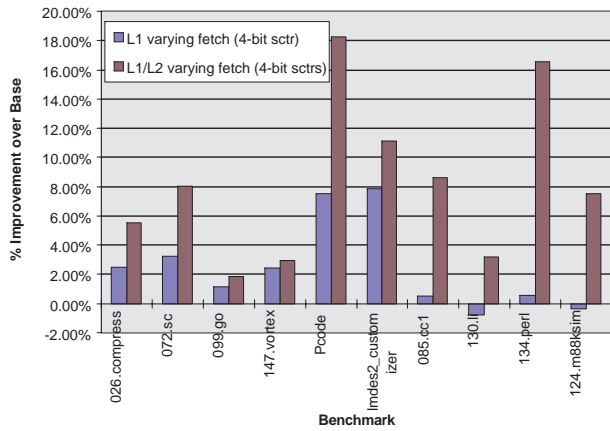


Figure 5: Performance for the spatial locality optimizations with 2-way and 4-way set-associative L1 and L2 data caches, respectively.

the optimizations at the L2 data cache level, due to the reductions we achieve in L2 cold start stall cycles for data with spatial locality.

### 5.3.3 Growing Memory Latency Effects

As discussed in Section 1, memory latencies are increasing, and this trend is expected to continue. Figure 6 shows the improvements achieved by our optimizations when applied to direct-mapped caches for both 100 and 200-cycle latencies, each relative to a base configuration with the same memory latency. Most of the benchmarks see much larger improvements from our optimizations, with the exception of *026.compress*. Because *026.compress* has very little spatial locality to exploit, the longer latency cannot be hidden as effectively. Although the raw number of cycles we eliminate grows, as a percentage of the associated base execution cycle count it becomes smaller.

### 5.3.4 Comparison of Integrated Techniques to Doubled Data Caches

As the memory latencies increase, intelligent cache management techniques will become increasingly important. We examined the performance improvement achieved by integrating our spatial locality optimizations with intelligent bypassing, using 8-bit access counters in each MAT entry [2]. The 4-way set-associative buffers used to hold the bypassed data at the L1 and L2 caches contain 128 8-byte entries and 512 32-byte entries, respectively. Then, the SLDT and MAT at each cache level are used to detect spatial locality and control the fetch sizes for both the data cache and the bypass buffer at that level.

Figure 7 shows the improvements achieved by combining these techniques at both cache levels for a 100-cycle memory latency. We show results for three direct-mapped MAT sizes: infinite, 1K-entry, and 512-entry. Also shown are the performance improvements achieved by doubling both the L1 and L2 data caches. Doubling the caches is a brute-force technique used to improve cache performance. Fig-

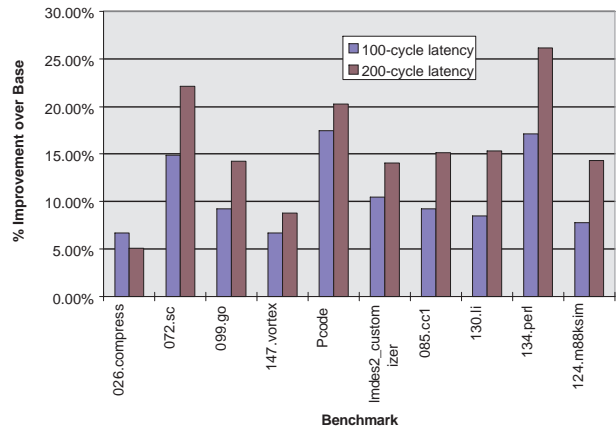


Figure 6: Performance for the spatial locality optimizations with growing memory latencies.

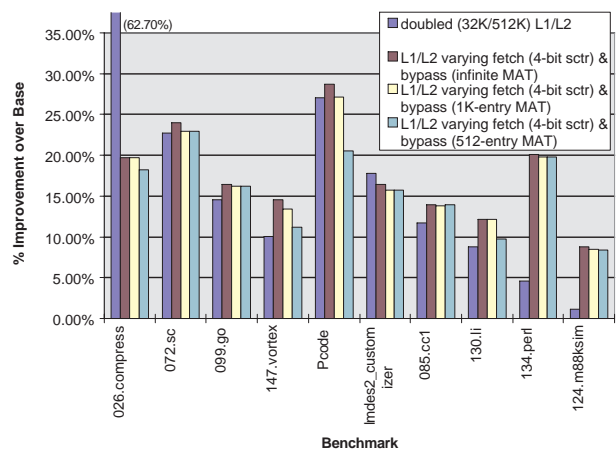


Figure 7: Comparison of doubled caches to integrated spatial locality and bypassing optimizations. Infinite, 1024-entry, and 512-entry direct-mapped MATs are examined.

ure 7 shows that performing our integrated optimizations at both cache levels can outperform simply doubling both levels of cache. The only case where the doubled caches perform significantly better than our optimizations is for *026.compress*. This improvement mostly comes from doubling the L2 data cache, which results because its hash tables can fit into a 512K-byte cache. *Pcode* is the only benchmark for which the performance degrades significantly when reducing the MAT size, however, 1K-entry MATs can still outperform the doubled caches. Comparing Figure 7 to the bypassing improvements in [2] shows that often significant improvements can be achieved by intelligently controlling the fetch sizes into the data caches and bypass buffers.

It can be shown that our optimizations require 26% and 44% less tags and data than doubling the data caches at the L1 and L2 levels, respectively [18]. The cost of our optimizations includes both the tag and data costs of the reorganized data caches (which have larger tag costs than the base configuration), the SLDTs, the MATs, and the bypass buffers. Comparing the performance of the spatial locality and by-

passing optimizations to the performance obtained by doubling the data caches at both levels, as shown in Figure 7, illustrates that for much smaller hardware costs our optimizations usually outperform simply doubling the caches.

## 6 Conclusion

Spatial locality optimizations must be able to detect and adapt to the varying amount of spatial locality both within and across applications in order to be effective. We presented a scheme which meets these objectives by detecting the amount of spatial locality in different portions of memory, and making dynamic decisions on the appropriate number of blocks to fetch on a memory access. A Spatial Locality Detection Table (SLDT), introduced in this paper, facilitates spatial locality detection for data while it is cached. This information is later recorded in a Memory Address Table (MAT) for long-term tracking, and is then used to tune the fetch sizes for each missing access.

Detailed simulations of several applications showed that significant speedups can be achieved by our techniques. The improvements are due to the reduction of conflict and capacity misses by utilizing small blocks and small fetch sizes when spatial locality is absent, and utilizing the prefetching effect of large fetch sizes when spatial locality exists. In addition, we showed that the speedups achieved by this scheme increase as the memory latency increases.

As memory latencies increase, the importance of cache performance improvements at each level of the memory hierarchy will continue to grow. Also, as the available chip area grows, it makes sense to spend more resources to allow intelligent control over the cache management, in order to adapt the caching decisions to the dynamic accessing behavior. We believe that our schemes can be extended into a more general framework for intelligent runtime management of the cache hierarchy.

## Acknowledgements

The authors would like to thank Mark Hill, Santosh Abraham and Wen-Hann Wang, as well as all the members of the IMPACT research group, for their comments and suggestions which helped improve the quality of this research. This research has been supported by the National Science Foundation (NSF) under grant CCR-9629948, Intel Corporation, Advanced Micro Devices, Hewlett-Packard, SUN Microsystems, NCR, and the National Aeronautics and Space Administration (NASA) under Contract NASA NAG 1-613 in cooperation with the Illinois Computer Laboratory for Aerospace Systems and Software (ICLASS).

## References

- [1] K. Boland and A. Dollas, "Predicting and precluding problems with memory latency," *IEEE Micro*, pp. 59–66, August 1994.
- [2] T. L. Johnson and W. W. Hwu, "Run-time adaptive cache hierarchy management via reference analysis," in *Proceedings of the 24th International Symposium on Computer Architecture*, pp. 315–326, June 1997.
- [3] S. Przybylski, "The performance impact of block sizes and fetch strategies," in *Proceedings of the 18th International Symposium on Computer Architecture*, pp. 160–169, June 1990.
- [4] A. J. Smith, "Line (block) size choice for cpu cache memories," *IEEE Transaction on Computers*, vol. C-36, pp. 1063–1075, 1987.
- [5] F. Dahlgren, M. Dubois, and P. Stenstrom, "Fixed and adaptive sequential prefetching in shared memory multiprocessors," in *Proceedings of the 1993 International Conference on Parallel Processing*, pp. 56–63, August 1993.
- [6] A. J. Smith, "Cache memories," *Computing Surveys*, vol. 14, no. 3, pp. 473–530, 1982.
- [7] N. P. Jouppi, "Improving direct-mapped cache performance by the addition of a small fully-associative cache and prefetch buffers," in *Proceedings of the 17th International Symposium on Computer Architecture*, pp. 364–373, June 1990.
- [8] J.-L. Baer and T.-F. Chen, "An effective on-chip preloading scheme to reduce data access penalty," in *Proceeding of Supercomputing '91*, pp. 176–186, Nov. 1991.
- [9] S. Mehrotra and L. Harrison, "Quantifying the performance potential of a data prefetch mechanism for pointer-intensive and numeric programs," Tech. Rep. 1458, CSR, Univ. of Illinois, November 1995.
- [10] J. W. C. Fu, J. H. Patel, and B. L. Janssens, "Stride directed prefetching in scalar processors," in *Proc. 25th Ann. Conference on Microprogramming and Microarchitectures*, Dec. 1992.
- [11] A. K. Porterfield, *Software Methods for Improvement of Cache Performance on Supercomputer Applications*. PhD thesis, Department of Computer Science, Rice University, Houston, TX, 1989.
- [12] T. C. Mowry, M. S. Lam, and A. Gupta, "Design and evaluation of a compiler algorithm for prefetching," in *Proc. Fifth Int'l Conf. on Architectural Support for Prog. Lang. and Operating Systems.*, pp. 62–73, Oct. 1992.
- [13] W. Y. Chen, S. A. Mahlke, P. P. Chang, and W. W. Hwu, "Data access microarchitectures for superscalar processors with compiler-assisted data prefetching," in *Proceedings of the 24th Annual International Symposium on Microarchitecture*, pp. 69–73, November 1991.
- [14] C.-K. Luk and T. C. Mowry, "Compiler-based prefetching for recursive data structures," in *Proceedings of the 7th International Conference on Architectural Support for Programming Languages and Operating Systems*, pp. 222–233, September 1996.
- [15] M. H. Lipasti, W. J. Schmidh, S. R. Kunkel, and R. R. Roediger, "SPAID: Software prefetching in pointer- and call-intensive environments," in *Proceedings of the 28th Annual International Symposium on Microarchitecture*, pp. 231–236, December 1995.
- [16] A. González, C. Aliagas, and M. Valero, "A data cache with multiple caching strategies tuned to different types of locality," in *Proc. International Conference on Supercomputing*, pp. 338–347, July 1995.
- [17] V. Milutinovic, B. Markovic, M. Tomasevic, and M. Tremblay, "The split temporal/spatial cache: Initial performance analysis," in *Proceedings of the SC'96*, March 1996.
- [18] T. L. Johnson, M. C. Merten, and W. W. Hwu, "Run-time spatial locality detection and optimization," Tech. Rep. IMPACT-97-02, University of Illinois, Urbana, IL (<http://www.crhc.uiuc.edu/IMPACT/papers/tech.html>), Sept. 1997.
- [19] R. A. Sugumar and S. G. Abraham, "Efficient simulation of caches under optimal replacement with applications to miss characterization," Tech. Rep. CSE-TR-143-92, Univ. of Michigan, 1992.



- [20] G. Tyson, M. Farrens, J. Matthews, and A. R. Pleszkun, "A modified approach to data cache management," in *Proceedings of the 28th Annual International Symposium on Microarchitecture*, pp. 93–103, December 1995.
- [21] J. A. Rivers and E. S. Davidson, "Reducing conflicts in direct-mapped caches with a temporality-based design," in *Proceedings of the 1996 International Conference on Parallel Processing*, pp. 151–162, August 1996.
- [22] J. W. C. Fu and J. H. Patel, "Data prefetching in multi-processor vector cache memories," in *Proc. 18th Ann. Int'l Symp. Computer Architecture*, pp. 54–63, June 1991.
- [23] P. P. Chang, S. A. Mahlke, W. Y. Chen, N. J. Warter, and W. W. Hwu, "IMPACT: An architectural framework for multiple-instruction-issue processors," in *Proceedings of the 18th International Symposium on Computer Architecture*, pp. 266–275, May 1991.
- [24] J. W. C. Fu and J. H. Patel, "How to simulate 100 billion references cheaply," Tech. Rep. CRHC-91-30, Center for Reliable and High-Performance Computing, University of Illinois, Urbana, IL, 1991.