

Exploiting the Prefetching Effect Provided by Executing Mispredicted Load Instructions

Resit Sendag *, David J. Lilja *, and Steven R. Kunkel **

* Department of Electrical and Computer Engineering
Minnesota Supercomputing Institute
University of Minnesota
200 Union St. S.E. Minneapolis, MN 55455
{rsgt, lilja}@ece.umn.edu

** IBM Corporation
3605 Highway 52 N.
Rochester, MN 55901
srkunkel@us.ibm.com

Abstract. *As the degree of instruction-level parallelism in superscalar architectures increases, the gap between processor and memory performance continues to grow requiring more aggressive techniques to increase the performance of the memory system. Several data prefetching techniques have been proposed for hiding the latency of main memory accesses, all of which must be implemented in such a way that prefetches are timely, useful and introduce little overhead. We propose a new technique, which is based on the wrong-path execution of loads far beyond instruction fetch-limiting conditional branches, to exploit more instruction-level parallelism by reducing the impact of memory delays. We examine the effects of the execution of loads down the wrong branch path on the performance of an aggressive issue processor. We find that, by continuing to execute the loads issued in the mispredicted path, even after the branch is resolved, we can actually reduce the cache misses observed on the correctly executed path. This wrong-path execution of loads can result in a speedup of up to 5% due to an indirect prefetching effect that brings data or instruction blocks into the cache for instructions subsequently issued on the correctly predicted path. However, it also can increase the amount of memory traffic and can pollute the cache. We propose the Wrong Path Cache (WPC) to eliminate the cache pollution caused by the execution of loads down mispredicted branch paths. For the configurations tested, fetching the results of wrong path loads into a fully associative 8-entry WPC can result in a 12% to 39% reduction in L1 data cache misses and in a speedup of up to 37%, with an average speedup of 9% over the baseline processor. Furthermore, the WPC consistently produces higher speedups than using a conventional victim cache [10].*

1 INTRODUCTION

Several methods have been proposed to exploit more instruction-level parallelism in superscalar processors and to hide the latency of the main memory accesses, including speculative execution [1-8] and data prefetching [9-26]. To achieve high issue rates, instructions must be fetched beyond the basic block-ending conditional branches. This can be done by speculatively executing instructions beyond branches until the branches are resolved. If the prediction was incorrect, the processor state must be restored to the state prior to the predicted branch and execution must be restarted down the correct path.

Processors with high instruction issue capabilities will speculatively execute many instructions to enable instruction issue to continue during the branch computation. This speculative execution will allow many memory references to be issued that turn out to be unnecessary since they are issued from the mispredicted branch path. However, these

incorrectly issued memory references may produce an indirect prefetching effect by bringing data or instruction lines into the cache that are needed later by instructions that are subsequently issued along correct execution path. On the other hand, these incorrectly issued memory references will increase the amount of memory traffic and can potentially pollute the cache with unneeded cache blocks [2].

Existing processors with deep pipelines and wide issue units do allow memory references to be issued speculatively down wrongly-predicted branch paths. In this study, however, we go one step further and examine the effects of continuing to execute the loads down the mispredicted branch path *even after the branch is resolved*. That is, we allow all speculatively issued loads to access the memory system if there is an available memory port. These instructions are marked as being from the mispredicted branch path when they are issued so they can be squashed in the write-back stage of the processor pipeline to prevent them from altering the target register after they access the memory system. In this manner, the processor is allowed to continue accessing memory with loads that are known to be from the wrong branch path. No store instructions are allowed to alter the memory system, however, since they are known to be invalid. The stores that are known to be down the wrong path after the branch is resolved are not executed eliminating the need for an additional speculative write buffer.

While this technique very aggressively issues load instructions to produce a significant impact on cache behavior, it has very little impact on the implementation of the processor's pipeline and control logic. The execution of wrong-path loads can make a significant performance improvement with very low overhead when there exists a large disparity between the processor cycle time and the memory speed. However, executing these loads can reduce performance in systems with small data caches and low associativities due to cache pollution. This cache pollution occurs when the wrong-path loads move blocks into the data cache that are never needed by the correct execution path. It also is possible for the cache blocks fetched by the wrong-path loads to evict blocks that still are required by the correct path.

In order to eliminate the cache pollution caused by the execution of the wrong-path loads, we propose the *Wrong Path Cache (WPC)*. This small fully-associative cache is accessed in parallel with the L1 cache. It buffers the values fetched by the wrong-path loads plus the castouts from the data cache. Our simulations show that the WPC can be very effective in eliminating the pollution misses caused by the execution of wrong path loads while simultaneously reducing the conflict misses that occur in the L1 data cache.

The remainder of the paper is organized as follows -- Section 2 describes the proposed wrong path cache. In Section 3, we present the details of the simulation environment with the simulation results given in Section 4. Section 5 discusses some related work. Section 6 suggests some future work with the conclusions given in Section 7.

2 WRONG PATH CACHE (WPC)

For small direct-mapped data caches, the execution of loads down the incorrectly-predicted branch path can reduce performance since the cache pollution caused by these wrong-path loads might offset the benefits of their indirect prefetching effect. Intuitively, we might expect that, the smaller the data cache, the higher the positive effect of the execution of loads down the wrong path since a larger number of cache misses will produce more indirect prefetches. However, some of these wrong-path loads will move some blocks into the data cache that are never needed by the correct execution path and will likely replace some of the blocks required by the correct path. This effect is likely to be more pronounced for low associativity caches.

In order to take advantage of the indirect prefetching effect of the wrong-path loads, we must eliminate the pollution they cause. We propose the Wrong Path Cache (WPC) specifically for this purpose. The idea is simply to use a small fully associative cache that is separate from the data cache to store the values returned by loads that are executed down the incorrectly-predicted branch path, as shown in Figure 1. Note that the WPC handles the loads that are *known* to be issued from the wrong path, that is, after the branch result is known. The loads that are executed before the branch is resolved are speculatively put in the L1 data cache.

If a wrong-path load causes a miss in the data cache, the required cache block is brought into the WPC instead of the data cache. The WPC is queried in parallel with the data cache when there is a miss in the data cache. The block is transferred simultaneously to the processor and the data cache when it is not in the data cache but it is in the WPC. When the address requested by a wrong-path load is in neither the data cache nor the WPC, the next cache level in the memory hierarchy is accessed. The required cache block is then placed only into the WPC to eliminate the pollution in the data cache that could otherwise be caused by the wrong-path loads. Note that misses due to loads on the correct execution path, and misses due to the loads issued from the wrong path before the branch is resolved, move the data into the data cache but not into the WPC.

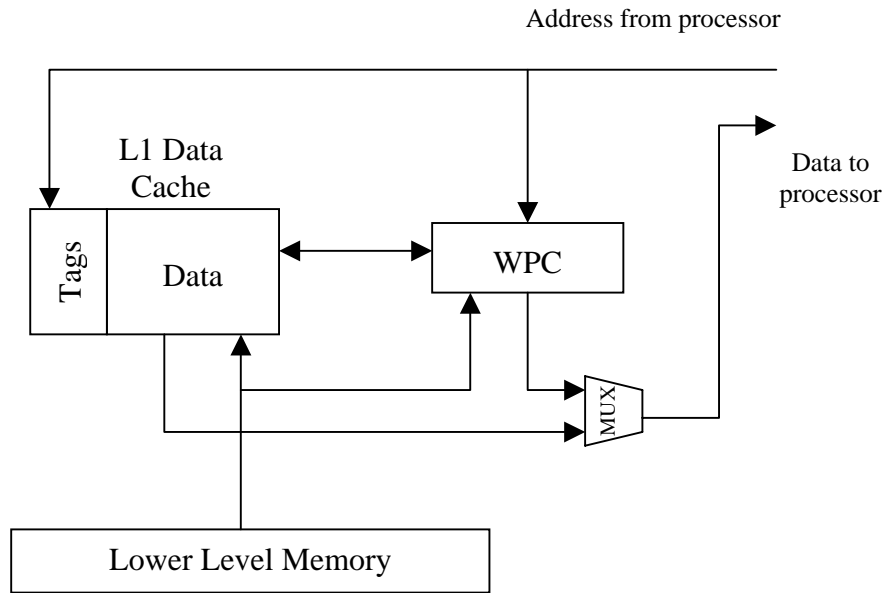


Figure 1. Placement of the WPC in the memory hierarchy.

The WPC also caches copies of blocks recently evicted by cache misses. That is, if the data cache must evict a block to make room for a newly referenced block, the evicted block is transferred to the WPC, as is done in the victim cache [10].

To summarize, the WPC is a small (4 to 16 entries for our experiments), fully associative cache that stores only the values returned by wrong-path loads and the values evicted from the cache. The basic algorithm for accessing the WPC is given in Figure 2. The design is essentially a combination of a cache for wrong-path loads and a victim cache [10].

```

If (wrong path execution)
    If (L1 data cache miss)
        If (WPC miss)
            Bring the block from the next level memory into the WPC;
        else // WPC hit
            NOP; // Update LRU info for the WPC
    else // L1 data cache hit
        NOP; // Update LRU info for the L1 data cache
else // WPC acts like a victim cache for this (correct) path
    If (L1 data cache miss)
        If (WPC miss)
            Bring the block from next level memory into L1 data cache
            Put the victim block into the WPC;
        else // WPC hit
            Swap the victim block and the WPC block;
    else // L1 hit
        NOP; // Update LRU info for the L1 data cache

```

Figure 2. The basic algorithm for accessing the WPC.

3 EXPERIMENTAL SETUP

3.1. Microarchitecture

Our microarchitectural simulator is built on top of the SimpleScalar toolset [27], version 3.0. SimpleScalar is an execution-driven simulator based upon the MIPS-IV ISA [28]. It supports the speculative execution of instructions, which is needed for this evaluation of the WPC. The simulator is modified to compare several different processor configurations, which are described in Section 3.2. The WPC is placed in the memory hierarchy as shown in Figure 1, and operates as described in Figure 2.

While fully-associative caches are expensive to build in terms of chip area, the small size of this supplemental cache makes it feasible to implement it on-chip, alongside the main L1 data cache. The access time of the WPC will be comparable to that of the much larger L1 cache. The multiplexer that selects between the WPC and the L1 cache could add a small delay to this access path, although this additional small delay also will occur with a victim cache.

The processor/memory model used in this study is an aggressively pipelined processor capable of issuing 8 instructions per cycle with out-of-order execution. It has a 128-entry reorder buffer with a 64-entry load/store buffer. The store forwarding latency is increased to 3 cycles in order to compensate for the added complexity of disambiguating loads and stores in a large execution window. There is a 3-cycle branch misprediction penalty. The processor has 8 integer ALU units, 2-integer MULT/DIV units, 4 load/store units, 6-FP Adders and 2-FP MULT/DIV units. The latencies are: ALU=1 cycle, MULT=3 cycles, integer DIV=12 cycles, FP Adder=2 cycles, FP MULT=4 cycles, and FP DIV=12 cycles. All the functional units, except the divide units, are fully pipelined to allow a new instruction to initiate execution each cycle.

The processor has a first-level 32 KB, 2-way set associative instruction cache. Various sizes of the L1 data cache (4KB, 8KB, 16KB, 32KB) with various associativities (direct-mapped, 2-way set-associative, 4-way set-associative) are examined in the following simulations. The first-level data cache is non-blocking with 4 ports. Both caches have block sizes of 32 bytes and 1-cycle hit latency. Since the memory footprints of the benchmark

programs used in this paper are somewhat small, a relatively small 256K 4-way associative unified L2 cache is used for all of the experiments in order to produce significant L2 cache activity. The L2 cache has 64-byte blocks and a hit latency of 12 cycles. The round-trip main memory access latency is 200 cycles for all of the experiments, unless otherwise specified. The effects on the WPC performance of varying the memory latency and the cache block size are examined in the simulations. There is a 64-entry 4-way set associative instruction TLB and 128-entry 4-way set associative data TLB, each with a 30-cycle miss penalty. For this study, we used the GAp branch predictor [30, 32]. The predictor has a 4K-entry Pattern History Table (PHT) with 2-bit saturating counters.

The WPC has 8-entries for all of the simulations, except for the experiments evaluating the effect of varying the size of the WPC.

3.2. Processor Configurations Tested

The following superscalar processor configurations are simulated to determine the performance impact of executing wrong-path loads, and the performance contributions of the Wrong Path Cache. The configurations, *all*, *vc*, and *wpc*, are modifications of the SimpleScalar [16] baseline processor described above.

orig: This configuration is the SimpleScalar baseline processor. It is an 8-issue processor with out-of-order execution and support for speculative execution of instructions issued from a predicted branch path. Note that this processor can execute loads from a mispredicted branch path. These loads can potentially change the contents of the cache, although they cannot change the contents of any registers. These wrong-path loads are allowed to access the cache memory system until the branch result is known. After the branch is resolved, they are immediately squashed and the processor state is restored to the state prior to the predicted branch. The execution then is restarted down the correct path. Wrong path loads that are waiting for their effective address to be computed, or are waiting for a free port to access the memory, before the branch is resolved do not access the cache and have no impact on the memory system.

all: In this configuration, the processor allows as many fetched loads as possible to access the memory system regardless of the predicted direction of conditional branches. This configuration is a good test of how the execution of the loads down the wrong branch path affects the memory system. Note that, in contrast to the *orig* configuration, the loads down the mispredicted branch direction are allowed to continue execution even after the branch is resolved. Wrong-path loads that are not ready to be issued before the branch is resolved, either because they are waiting for the effective address calculation or for an available memory port, are issued to the memory system if they become ready after the branch is resolved, even though they are known to be from the wrong path. Instead of being squashed after the branch is resolved as in the *orig* configuration, they are allowed to access the memory. However, they are squashed before being allowed to write to the destination register. Note that a wrong-path load that is dependent upon another instruction that gets flushed after the branch is resolved also is flushed in the same cycle. Wrong-path stores are not allowed to execute in this configuration to eliminate the need for an additional speculative write buffer. Stores are squashed as soon as the branch result is known.

orig_vc: This configuration is the *orig* configuration (the baseline processor) with the addition of an 8-entry victim cache.

all_vc: This configuration is the *all* configuration (the baseline processor) with the addition of an 8-entry victim cache. It is used to compare against the performance improvement made possible by caching of the wrong-path loads in the WPC.

wpc: This configuration adds an 8-entry Wrong Path Cache (WPC) to the *all* configuration.

3.3. Benchmark Programs

The test suite used in this study consists of the combination of SPEC95 and SPEC2000 benchmark programs shown in Table 1. All benchmarks were compiled using gcc 2.6.3 at optimization level O3 and each benchmark ran to completion.

The SPEC2000 benchmarks are run with specially developed input data sets that limit their total simulation time while maintaining the fundamental characteristics of the programs' overall behaviors. For more information on the reduced input sets for these benchmarks, see [29].

TABLE 1. Selected benchmark characteristics for the benchmarks used in this study.

Benchmark	Suite	Type	Instructions(M)	Input set
<i>124.m88ksim</i>	SPEC95	Integer	120.1	Train
<i>126.gcc</i>	SPEC95	Integer	873.2	Test
<i>130.li</i>	SPEC95	Integer	183.3	Train
<i>132.ijpeg</i>	SPEC95	Integer	553.3	Test
<i>134.perl</i>	SPEC95	Integer	2191.4	Test
<i>164.gzip</i>	SPEC2000	Integer	763.6	Reduced Medium
<i>175.vpr</i>	SPEC2000	Integer	216.9	Reduced Medium
<i>181.mcf</i>	SPEC2000	Integer	202.1	Reduced Medium
<i>197.parser</i>	SPEC2000	Integer	513.7	Reduced Medium
<i>300.twolf</i>	SPEC2000	Integer	214.6	Reduced Medium
<i>183.quake</i>	SPEC2000	Floating-Point	715.9	Reduced Large

4 RESULTS

The simulation results are presented as follows. First, the performances of the different configurations are compared using the speedups relative to the baseline (*orig*) processor. Next, several important memory system parameters are varied to determine the sensitivity of the WPC to these parameters. The impact of executing wrong-path loads both with and without the WPC also is analyzed.

In this paper, our focus is on improving the performance of on-chip direct-mapped data caches. Therefore, most of the comparisons for the WPC are made against a victim cache [10]. We do investigate the impact of varying the L1 associativity in Section 4.2.5, however.

4.1 Performance Comparisons

4.1.1 Speedup Due to the WPC

Figure 3 shows the speedups obtained relative to the *orig* configuration when executing each benchmark on the different configurations described in Section 3.2. The results are given for an 8KB direct-mapped data cache using 32-byte cache blocks. The WPC and the victim cache each have eight entries in those configurations that include these structures.

Of all of the configurations, *wpc*, which executes loads down the wrong branch path with an 8-entry WPC, gives the greatest speedup. From Figure 3, we can see that, for small caches, the *all* configuration actually produces a slowdown due to the large number of wrong-path loads polluting the L1 cache. However, by adding the WPC, the new configuration, *wpc*, produces the best speedup compared to the other configurations. In particular, *wpc* outperforms the *orig_vc* and *all_vc* configurations, which use a simple victim cache to improve the performance of the baseline processor. While *wpc* can produce a speedup of up to 37% for the *mcf* benchmark compared to the baseline processor, it also produces consistently higher speedups compared to the configurations with a victim cache.

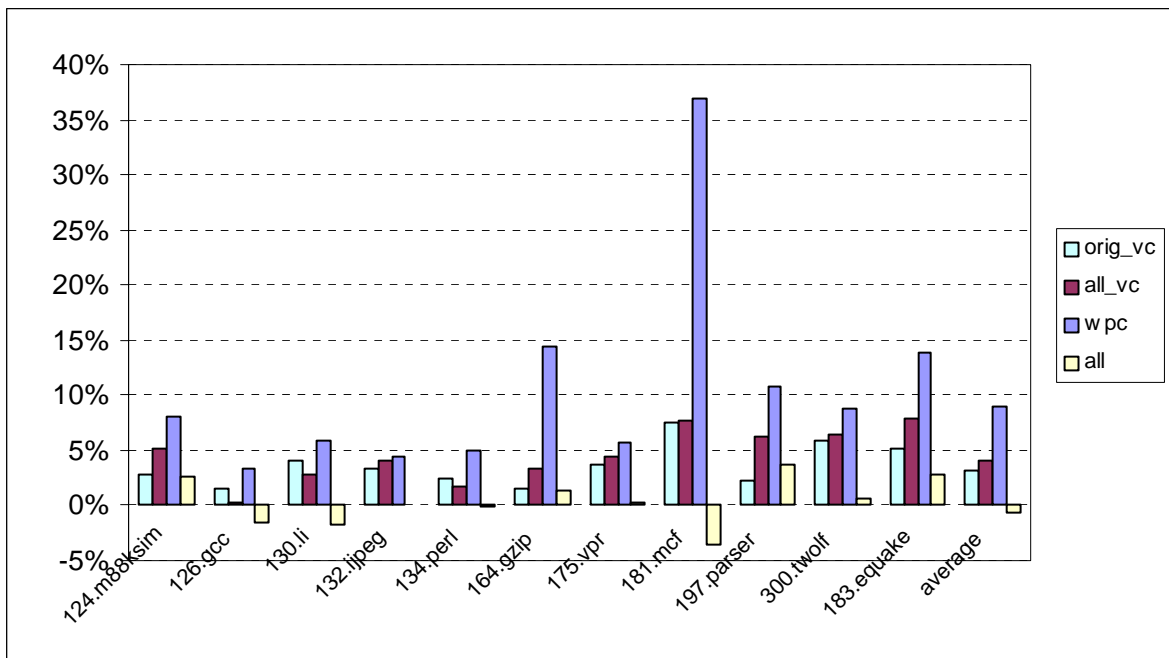


Figure 3. The Wrong Path Cache (*wpc*) produces consistently higher speedups than the victim cache (*vc*) or the *all* configuration, which does not have a WPC but does execute all ready wrong-path loads if there is a free port to the memory system. The data cache is 8KB direct-mapped and has 32 byte blocks. All speedups are relative to the baseline (*orig*) processor.

Figure 3 shows that the WPC produces better performance than a simple victim cache of the same size. With an 8KB direct-mapped data cache, the average speedup obtained from using the WPC (*wpc*) is substantially better than that obtained from using only a victim cache (*orig_vc* and *all_vc*) compared to the baseline processor (*orig*). While both the WPC and the victim cache reduce the impact of conflict misses in the data cache by storing recent castouts near the processor, the WPC goes further by acting like a prefetch buffer and thus preventing pollution misses due to the indirect prefetches caused by executing the wrong-path loads in the *all* configuration. The WPC also may unintentionally be reducing the latency of retrieving data

from other levels in the memory hierarchy for both compulsory and capacity misses from loads executed on the correct path. However, analyzing these small distinctions is beyond the scope of this paper.

While we will study the effect of different cache parameters in later sections, Figure 4 shows the speedup results for an 8KB L1 data cache with 4-way associativity. When increasing the associativity of the L1 cache, the speedup obtained by the *orig_vc* seen in Figure 3 disappears. However, the *wpc* still provides significant speedup as the associativity increases and it substantially outperforms the *all_vc* configuration. The *mcf* program shows generally poor cache behavior and increasing the L1 associativity does not reduce its miss rate significantly. Therefore, we see that the speedup produced by the *wpc* for *mcf* remains the same in Figures 3 and 4. As expected, a better cache with lower miss rates reduces the benefit of the *wpc*. From Figure 4, we also see that the *all* configuration can produce some speedup. There is still some slowdown for a few of the benchmarks due to pollution from the wrong path execution of loads. However, the slowdown for the *all* configuration is less than in Figure 3, where the cache is direct-mapped.

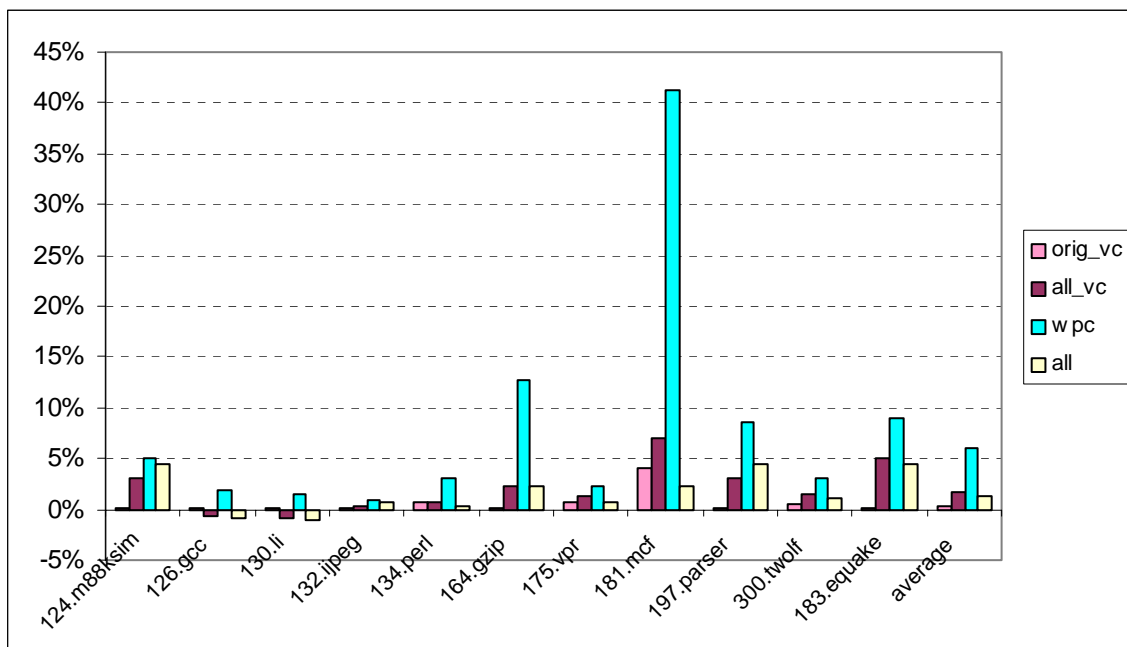


Figure 4. With a data cache of 8KB with 4-way associativity, the speedup obtained by *orig_vc* disappears. However, *wpc* continues to provide significant speedup and substantially outperforms the *all_vc* configuration. The *all* configuration also shows significant speedup for some benchmarks. The data cache has 32 byte blocks. All speedups are relative to the baseline (*orig*) processor.

4.1.2 Causes of the WPC Speedups

Figures 5-7 analyze how the execution of loads down the wrong path actually affects the memory system and how the WPC helps improve the overall performance. Figures 5 and 6 show how both the correct path and wrong path accesses to the memory hierarchy are serviced. In these figures, the *wrong path* is the execution path that is known to be wrong after the branch resolves. The *correct path*, on the other hand, includes the speculative execution of loads that are issued before the branch is resolved.

The speedup results shown in Figures 3 and 4 can be at least partially explained by examining what levels of the memory hierarchy service the memory accesses. Figure 5 shows

that the great majority of all memory accesses in the benchmark programs are serviced by the L1 cache, as is to be expected. While a relatively small fraction of the memory accesses cause misses, these misses add a disproportionately large amount of time to the memory access time. The values for memory accesses that miss in the L1 cache must be obtained from one of three possible sources, the wrong-path cache (WPC), the L2 cache, or the memory. Figure 5 shows that a substantial fraction of the misses in these benchmark programs are serviced by the WPC. For example, 4% of all memory accesses issued by *twolf* are serviced by the WPC. However, this fraction corresponds to 32% of the L1 misses generated by this program. Similarly, 3.3% of *mcf*'s memory accesses, and 1.9% of *equake*'s, are serviced by the WPC, which corresponds to 21% and 29% of their L1 misses, respectively. Since the WPC is accessed in parallel with the L1 cache, misses serviced by the WPC are serviced in the same amount of time as a hit in the L1 cache, while accesses serviced by the L2 cache require 12 cycles and accesses that must go all the way to memory require 200 cycles. For most of these programs, we see that the WPC converts approximately 20-35% of misses that would have been serviced by the L2 cache or the memory into accesses that are equivalent to an L1 hit.

While the above discussion explains some of the speedups seen in Figures 3 and 4, it does not completely explain the results. For instance, *twolf* has the largest fraction of memory accesses serviced by the WPC in Figure 5. However, *mcf*, *gzip*, and *equake* show better overall speedups. This difference in speedup is explained in Figure 6. This figure shows which levels of the memory hierarchy service the speculative loads issued on what is subsequently determined to be the wrong branch path. Speculative loads that miss in both the L1 cache and the WPC are serviced either by the L2 cache or by the memory. These values are placed in the WPC in the hope that the values will be subsequently referenced by a load issued on the correct branch path.

In Figure 6, we see that 30 percent of the wrong path accesses that miss in both the L1 and the WPC are serviced by memory, which means that this percentage of the blocks in the WPC are loaded from memory. So, from Figure 5 we can say that 30 percent of the correct path accesses that hit in the WPC for *mcf* would have been serviced by the memory in a system without the WPC. That is, the WPC effectively converts a large fraction of this program's L1 misses into the equivalent of an L1 hit. In *twolf*, on the other hand, most of the hits to the WPC would have been hits in the L2 cache in the absence of the WPC. We see in Figure 6 that less than 1% of the wrong path accesses for *twolf* that miss both in the L1 and the WPC are serviced by memory, while 99% of these misses are serviced by the L2 cache. That is, almost all the data in the WPC comes from the L2 cache for *twolf*. Thus, the WPC does a better job of hiding miss delays for *mcf* than for *twolf*, which explains why *mcf* obtains a higher overall speedup with the WPC than does *twolf*. A similar argument explains the speedup results observed in the rest of the programs, as well.

In Figure 7, we see that approximately 55-75% of the blocks that are moved into the WPC are used before being evicted by subsequent accesses. However, not all of these accesses are useful to the correct execution path. The lighter portion of each bar in this figure shows the fraction of all of the blocks accessed in the WPC that are actually used by loads issued from the correct execution path. For *jpeg*, for example, 55% of the blocks in the WPC are accessed before being evicted. Furthermore, of these blocks, only 19% are actually useful to the execution of the correct branch path. For *mcf*, on the other hand, 74% of the blocks in the WPC are accessed before being evicted, and 61% of these accessed blocks are useful to the correct path execution. *M88ksim* has a usage percentage of 64%, but at 6%, it has the smallest percentage used by the correct path. However, *m88ksim* still has a greater benefit from the WPC than *jpeg* because the overall percentage of correct path accesses serviced by the WPC is larger in *m88ksim* than *jpeg* (see Figure 5). In addition, *jpeg* has a smaller percentage of

wrong path references serviced by the L2 cache or memory, which implies that there is less prefetching for the later correct path accesses for *jpeg* compared to *m8ksim*.

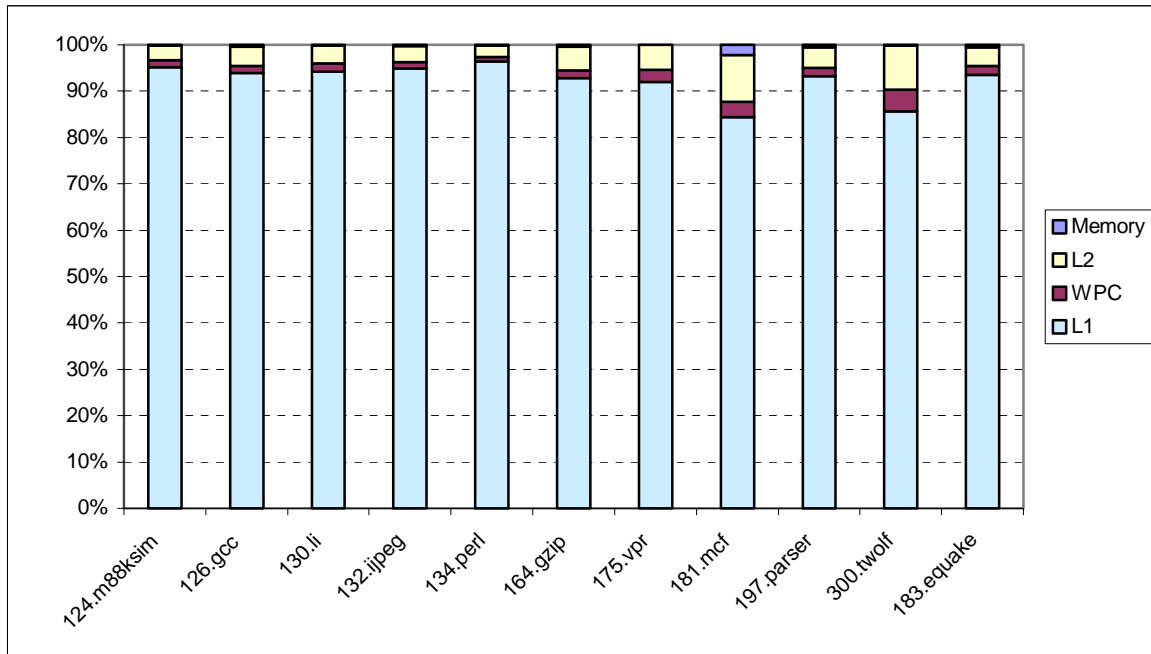


Figure 5. The fraction of memory references on the *correct* execution path that are serviced by the L1 cache, the WPC, the L2 cache, and memory. The L1 data cache is 8KB direct-mapped and has 32 byte blocks.

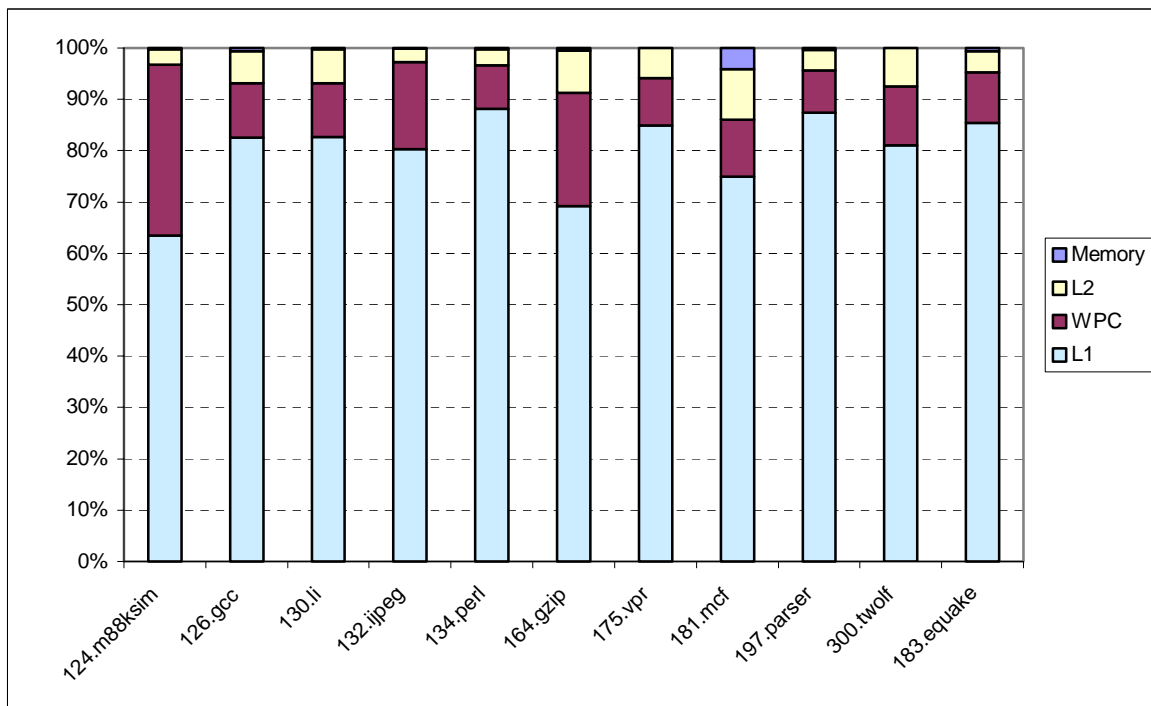


Figure 6. The fraction of memory references on the *wrong* execution path that are serviced by the L1 cache, the WPC, the L2 cache, and memory. The L1 data cache is 8KB direct-mapped and has 32 byte blocks.

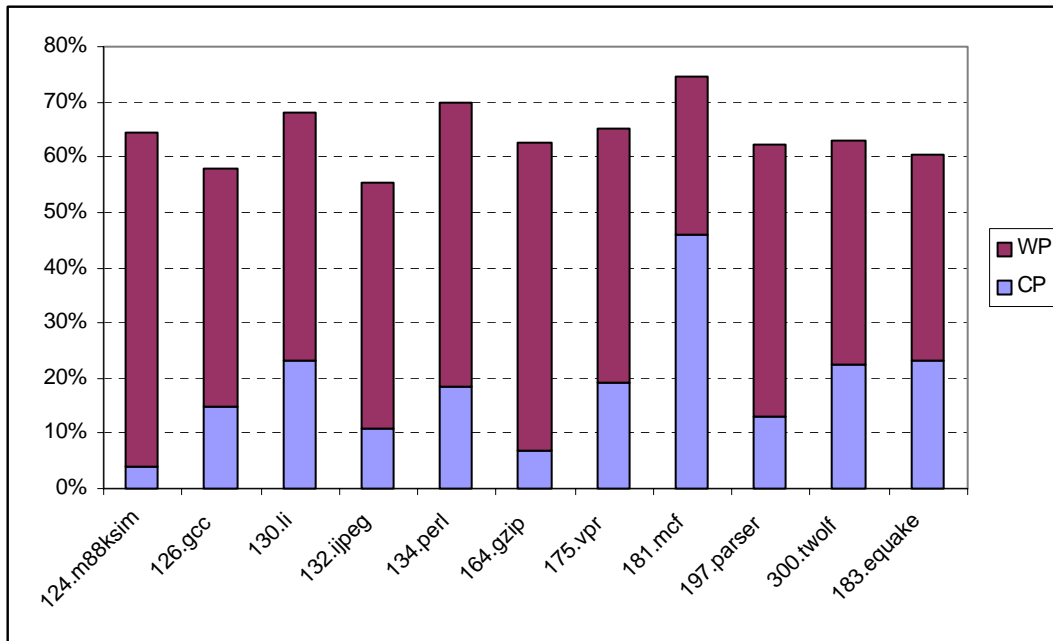


Figure 7. Percentage of the blocks that are moved into the WPC that are used before being evicted. The proportion of hits used by the Wrong Path (WP) and the Correct Path (CP) are also given.

4.2 Sensitivity to Cache Parameters

There are several parameters that affect the performance of a cache memory system. In this study, we examine the effects of the cache size, the associativity, the cache block size, and the memory latency on the cache performance when allowing the execution of wrong-path loads both with and without the WPC.

4.2.1. The Effect of the L1 cache size

Figure 8 shows that relative benefit of the *wpc* decreases as the L1 cache size increases. However, the WPC size is kept constant in these simulations so that the relative size of the WPC to the data cache is reduced. With a smaller cache, wrong-path loads cause more misses compared to configurations with larger caches. These additional misses tend to prefetch data that is put into the WPC for use by subsequently executed correct branch paths. The WPC eliminates the pollution in the L1 data cache for the *all* configuration that would otherwise have occurred without the WPC, which then makes these indirect prefetches useful for the correct branch path execution.

While the WPC is a relatively small hardware structure, it does consume some chip area. Figure 9 shows the performance obtained with an 8-entry WPC used in conjunction with an 8KB L1 cache compared to the performance obtained with the original processor configuration using a 16KB L1 cache and a 32KB L1 cache but without a WPC. We find that, for all of the test programs, the small WPC with the 8KB cache exceeds the performance of the processor when the cache size is doubled, but without the WPC. Furthermore, the WPC configuration exceeds the performance obtained when the size of the L1 cache is quadrupled for all of the test programs except *gcc*, *li*, *vpr*, and *twolf*. We conclude that this small WPC is an excellent use of the chip area compared to simply increasing the L1 cache size.

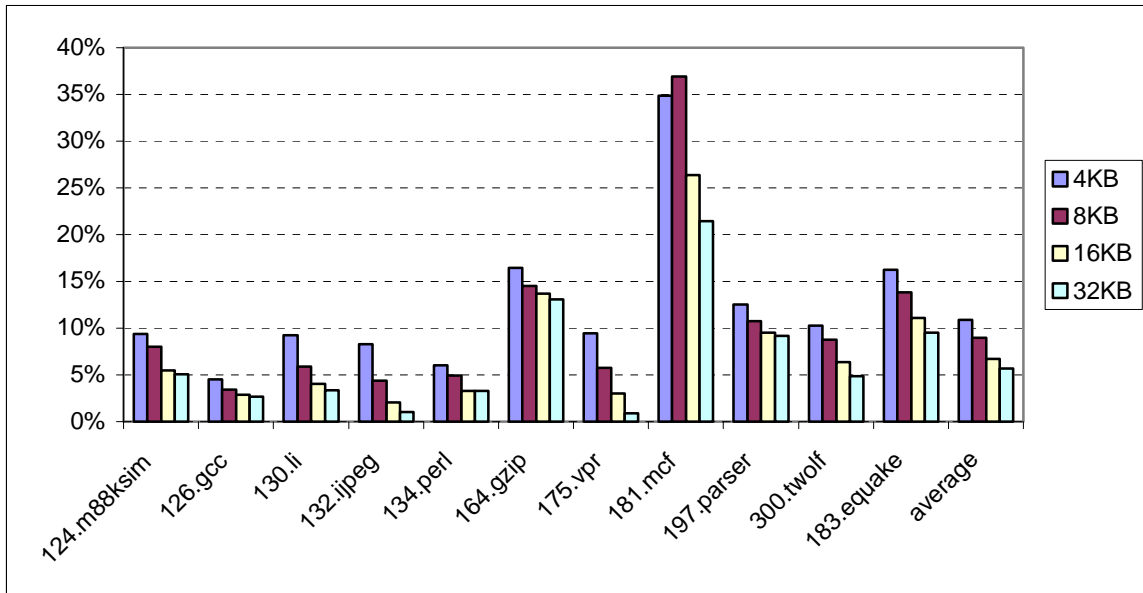


Figure 8. Speedup obtained with the *wpc* configuration as the L1 cache size is varied. The L1 data cache is direct-mapped with 32 byte blocks. All speedups are relative to the baseline (*orig*) processor.

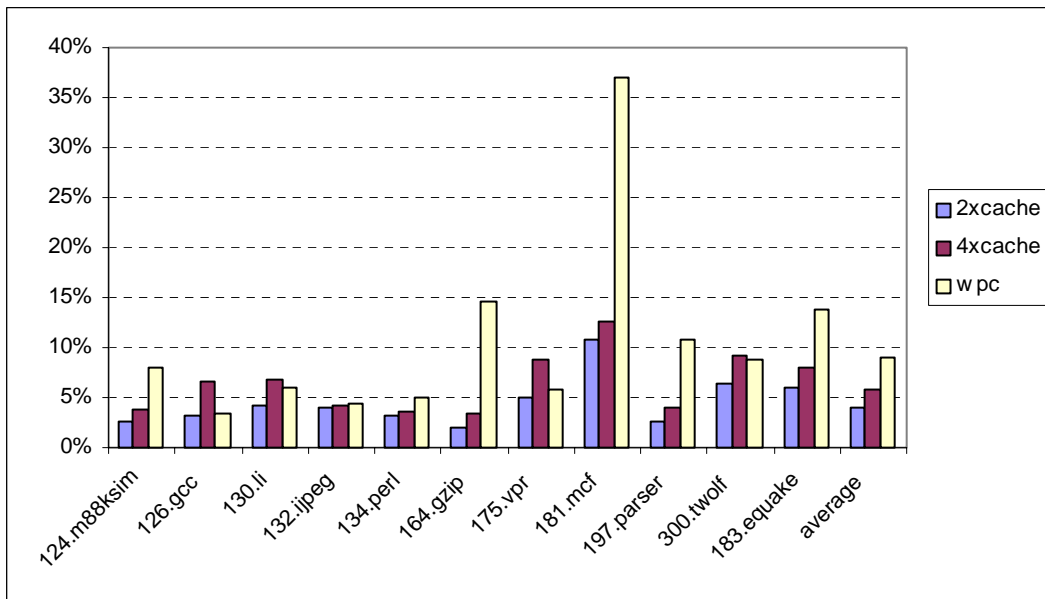


Figure 9. The speedup obtained with the WPC compared to configurations with larger L1 caches but without a WPC. The base cache size is 8KB and is direct-mapped with 32 byte blocks.

4.2.2. The Effect of the WPC size

The effect of varying the WPC size is shown in Figures 10 and 11 for WPC sizes of 4, 8, and 16 entries (128B, 256B, and 512B, respectively). The size of a straight victim cache also is varied to compare the victim cache to the WPC. Figure 10 shows the speedup for a direct-mapped L1 data cache, while Figure 11 shows the impact of varying sizes of the WPC for a 4-way associative L1 data cache. We see that even a small 4-entry WPC generally produces better performance than a 16-entry victim cache. In Figure 11, we can see that, for most of the

benchmark programs, there is no (or negligible) speedup with a straight victim cache (*orig_vc* and *all_vc*), but the *wpc* continues to show significant speedups.

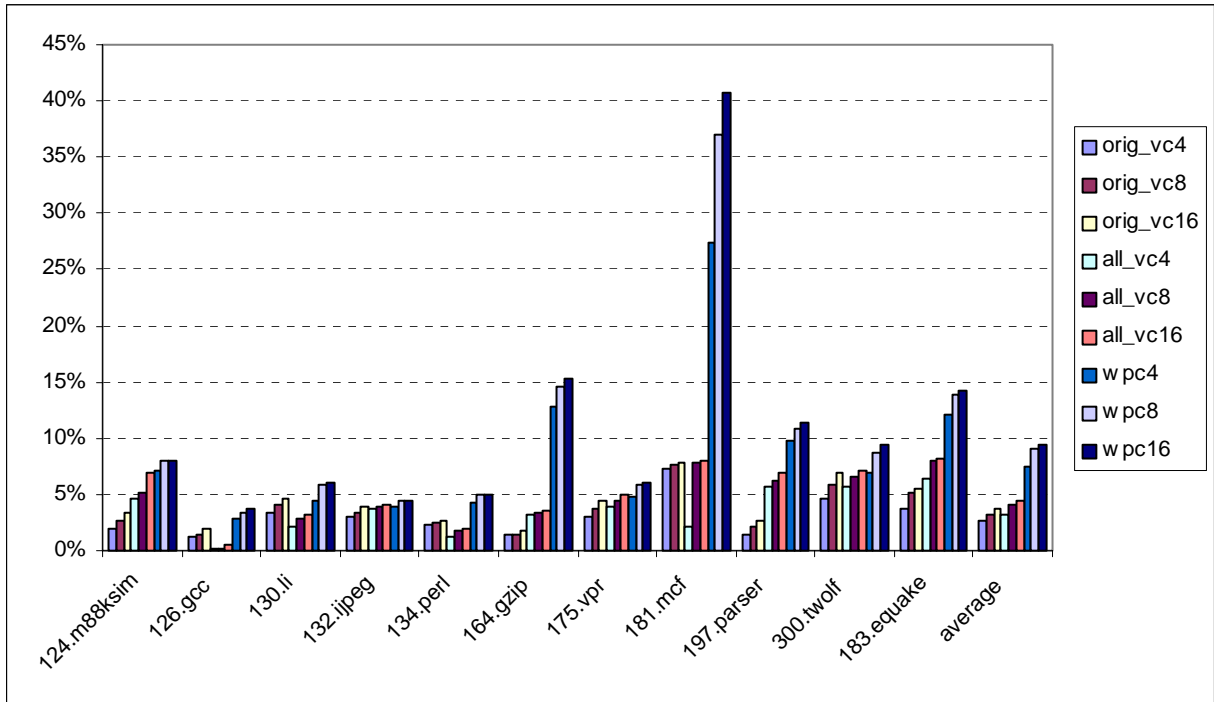


Figure 10. Speedups relative to the baseline (*orig*) processor when varying the size of the victim cache and the size of the WPC. The L1 data cache is 8 KB direct-mapped with 32 byte blocks. The legend denotes the number of entries in the victim cache or WPC for each design.

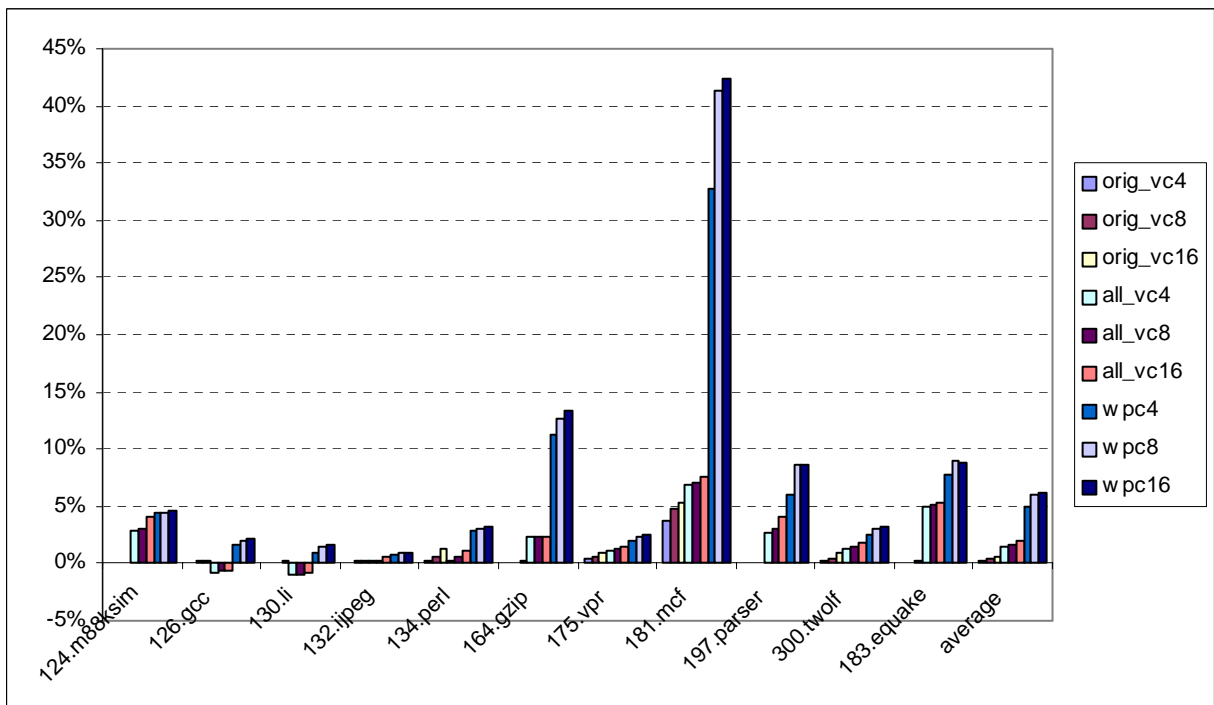


Figure 11. Speedups relative to the baseline (*orig*) processor when varying the size of the victim cache and the size of the WPC. The L1 data cache is 8 KB, 4-way associative with 32 byte blocks. The legend denotes the number of entries in the victim cache or WPC for each design.

4.2.3. The Effect of Memory Latency

In a traditional hardware- or software-based prefetching implementation [9-24], the target addresses must be fetched as part of the main execution path. However, since the prefetched value typically is needed almost immediately by an instruction on this execution path, there often is not enough time to cover the memory latency for the prefetched value. The execution of the wrong-path loads, on the other hand, indirectly prefetches down a path that is not immediately taken. As a result, these wrong-path loads potentially have more time to prefetch a block from memory before the correct path that actually needs the indirectly prefetched values is executed. Figure 12 shows that the WPC can effectively hide even relatively large main memory latencies due to this indirect prefetching effect.

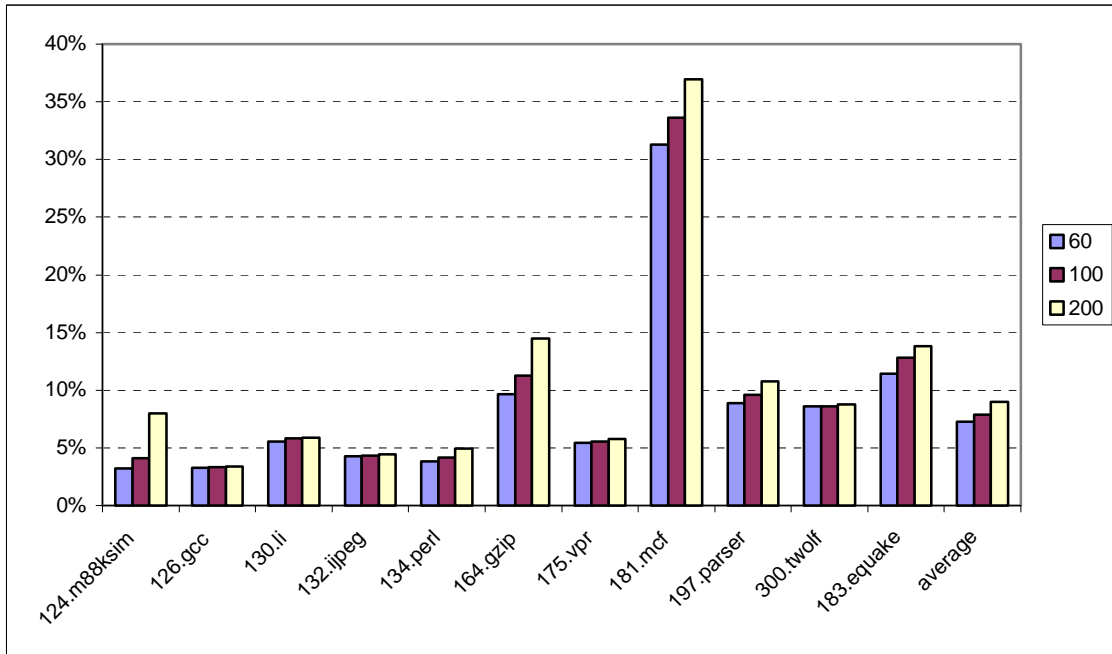


Figure 12. Speedup of the *wpc* configuration compared to the baseline (*orig*) configuration when varying the memory latency. The data cache is 8 KB direct-mapped with 32-byte cache blocks.

4.2.4. Total Data and Bus Traffic

Figure 13 shows that executing the loads that are known to be down the wrong path typically increases the number of L1 data cache references by about 15-25% for most of the test programs. Furthermore, this figure shows that executing these wrong-path loads increases the bus traffic (measured in bytes) between the L1 cache and the L2 cache by 5-23%, with an average increase of 11%. However, the WPC reduces the total data cache miss ratio for loads on the correct path by up to 39%, as shown in Figure 14. This reduction occurs because of the indirect prefetching effect of the wrong-path misses, which subsequently reduce the number of correct-path misses.

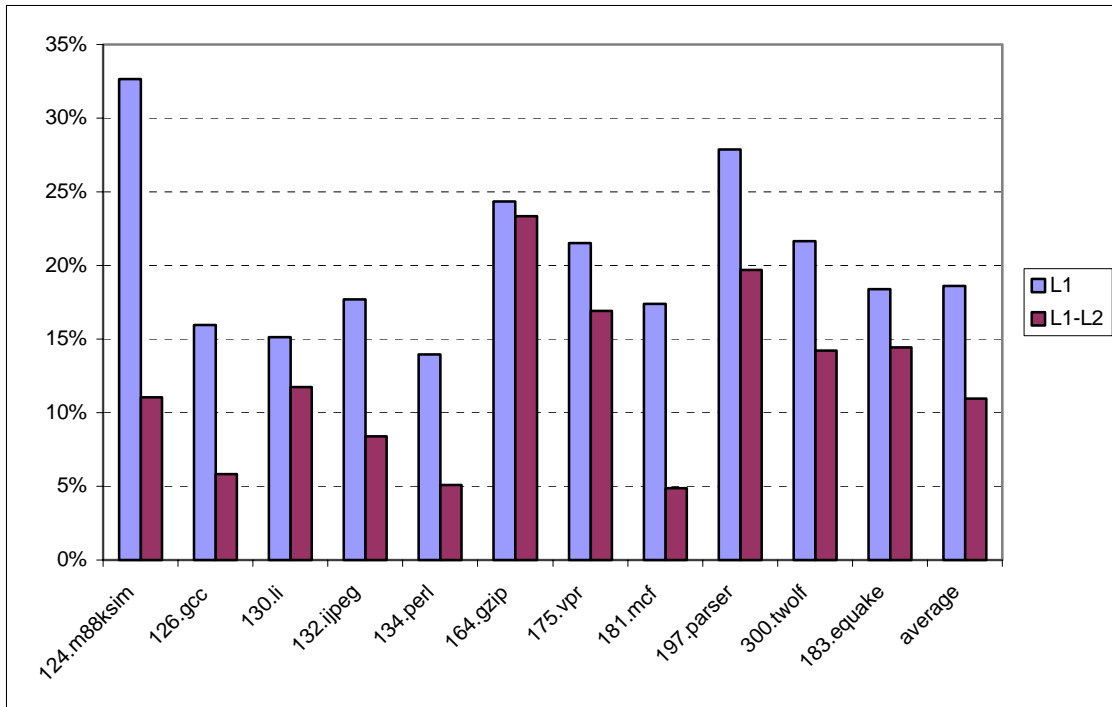


Figure 13. The percentage increase in L1 cache accesses and traffic between the L1 cache and the L2 cache for the *wpc* configuration compared to the *orig* configuration. The L1 cache is 8 KB, direct-mapped and has 32B blocks.

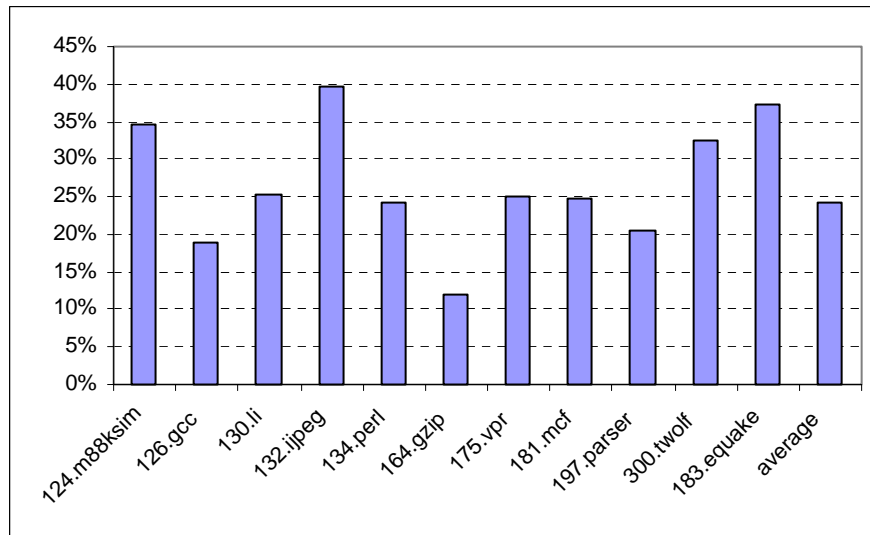


Figure 14. The reduction in data cache misses for the *wpc* configuration compared to the *orig* configuration. The L1 cache is 8 KB, direct-mapped and has 32B blocks.

4.2.5. Impact of L1 Cache Associativity

Increasing the L1 cache associativity typically tends to reduce the number of L1 misses on both the correct path [9] and the wrong path. This reduction in misses reduces the number of indirect prefetches issued from the wrong path, which then reduces the impact of the WPC, as shown in Figure 15. The *mcf* program is the exception since its overall cache behavior is less sensitive to the L1 associativity than the other test programs.

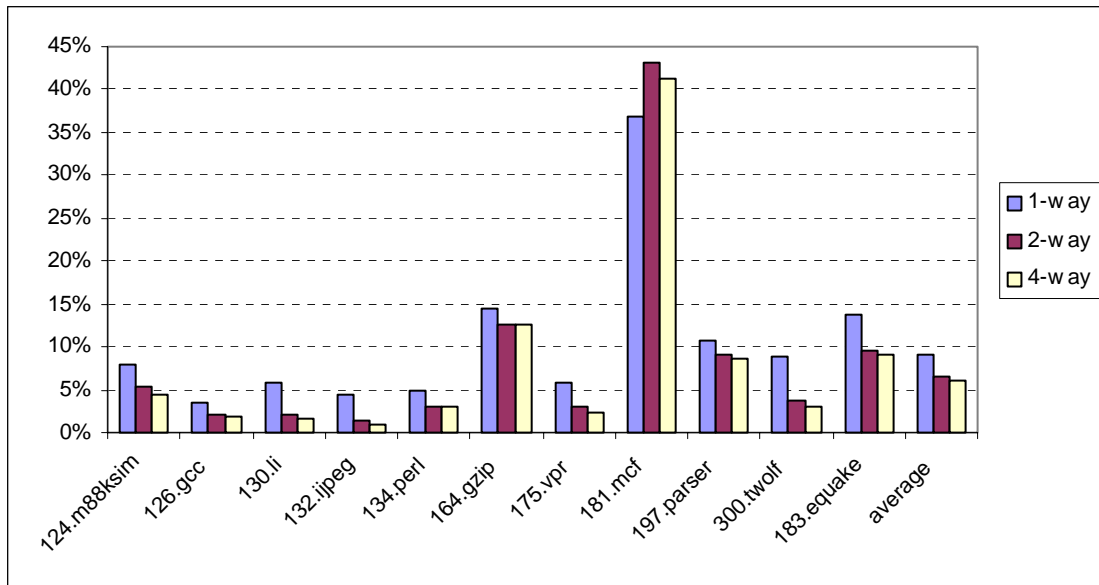


Figure 15. The effect of the L1 cache associativity on the speedup of the *wpc* configuration compared to the *orig* configuration. The L1 cache size is 8 KB with 32 byte blocks.

4.2.6. The Effect of L1 Cache Block Size

As the block size of the data cache increases, the number of conflict misses also tends to increase [9, 31]. Figure 16 shows that smaller cache blocks produce better speedups for configurations without a WPC when wrong-path loads are allowed to execute since larger blocks more often displace useful data in the L1 cache. However, for the systems with a WPC, the increasing conflict misses in the data cache due to the larger blocks results in an increase in these misses hitting in the WPC because of the victim-caching behavior of the WPC. In addition, the indirect prefetches provide a greater benefit for large blocks since the WPC eliminates their polluting effects. We conclude that larger cache blocks work well with the WPC since the strengths and weaknesses of larger blocks and the WPC are complementary.

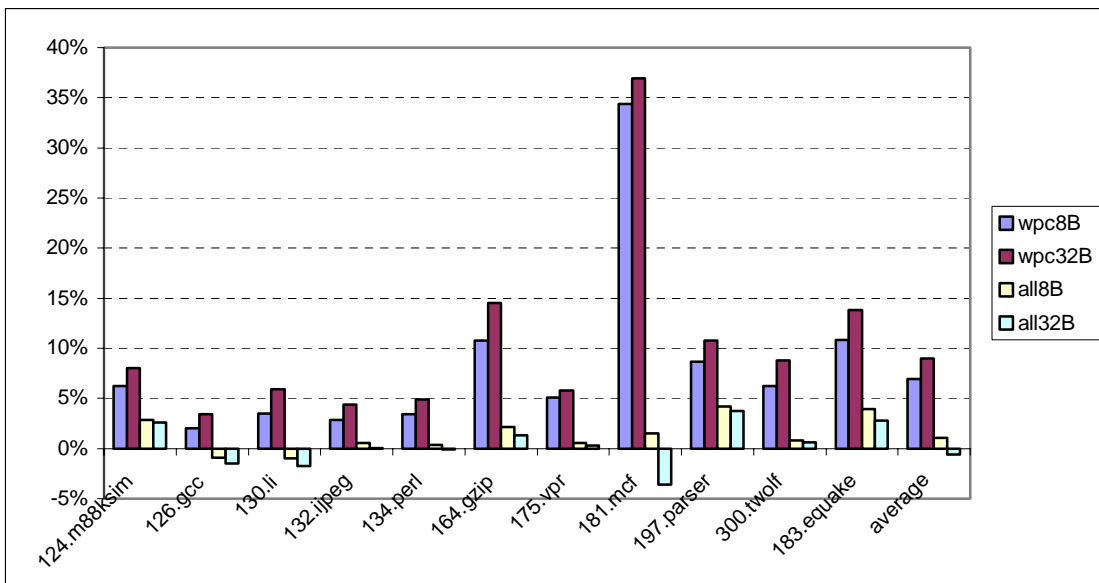


Figure 16. The effect of the cache block size on the speedup of the *all* and *wpc* configurations compared to the *orig* configuration. The L1 cache is direct-mapped and 8 KB. The WPC is 256B, i.e., 8-entries with 32B lines (*wpc32B*), or 32-entries with 8B lines (*wpc8B*).

4.3. The Impact of Branch Prediction Accuracy

Since the indirect prefetching effect provided by the WPC relies on the execution of load instructions issued along a wrongly-predicted branch path, it seems likely that the speedup due to the WPC would be sensitive to the accuracy of the branch predictor. For instance, higher branch prediction accuracies would result in fewer wrong-path loads being issued, which could reduce the benefit of the WPC.

To test this sensitivity, we examined the performance obtained with the WPC for three different branch predictors. The prediction accuracies of the different predictors are shown in Table 2. The *bpred2* predictor is the GAp branch predictor [30, 32] with a 4K-entry Pattern History Table (PHT) with 2-bit saturating counters, which has been used in all of the simulations up to this point. The *bpred1* predictor, which has the lowest prediction accuracy, uses a smaller PHT with 1K-entries. The *bpred3* predictor, which uses an 8K-entry PHT, has the highest prediction accuracy.

Table 2. Branch prediction accuracies of the three different branch predictor configurations examined in this study

Benchmark	%accuracy <i>bpred1</i>	%accuracy <i>bpred2</i>	%accuracy <i>bpred3</i>
<i>124.m88ksim</i>	86.8	93.0	95.6
<i>126.gcc</i>	76.9	80.5	85.6
<i>130.li</i>	87.7	90.3	92.4
<i>132.jpeg</i>	81.8	85.3	91.2
<i>134.perl</i>	84.0	89.7	94.4
<i>164.zip</i>	86.5	89.1	91.9
<i>175.vpr</i>	76.9	78.1	84.8
<i>181.mcf</i>	84.6	87.3	91.9
<i>197.parser</i>	85.7	88.6	92.8
<i>300.twolf</i>	71.4	74.4	80.1
<i>183.earthquake</i>	81.9	85.0	90.5

The speedup results in Figure 17 show that the performance obtained with the WPC does not directly correlate with the branch prediction accuracy. For half of the test programs, increasing the prediction accuracy reduces the speedup provided by the WPC. This reduction is presumably due to fewer misses occurring on the wrong branch path, which reduces the indirect prefetching effect provided by the WPC. For the remaining programs, however, increasing the prediction accuracy causes an increase in the speedup provided by the WPC. This result is counter-intuitive since there are fewer wrong-path loads issued as the branch prediction accuracy increases. It is not entirely clear why fewer wrong-path loads would cause an increase in the performance provided by the WPC. However, it appears likely that, for these test programs, many of the wrong-path loads issued for the configurations with lower prediction accuracies are not useful to the correct-path execution. Additional work still needs to be done to fully understand the relationship of the branch prediction accuracy to the performance of the WPC.

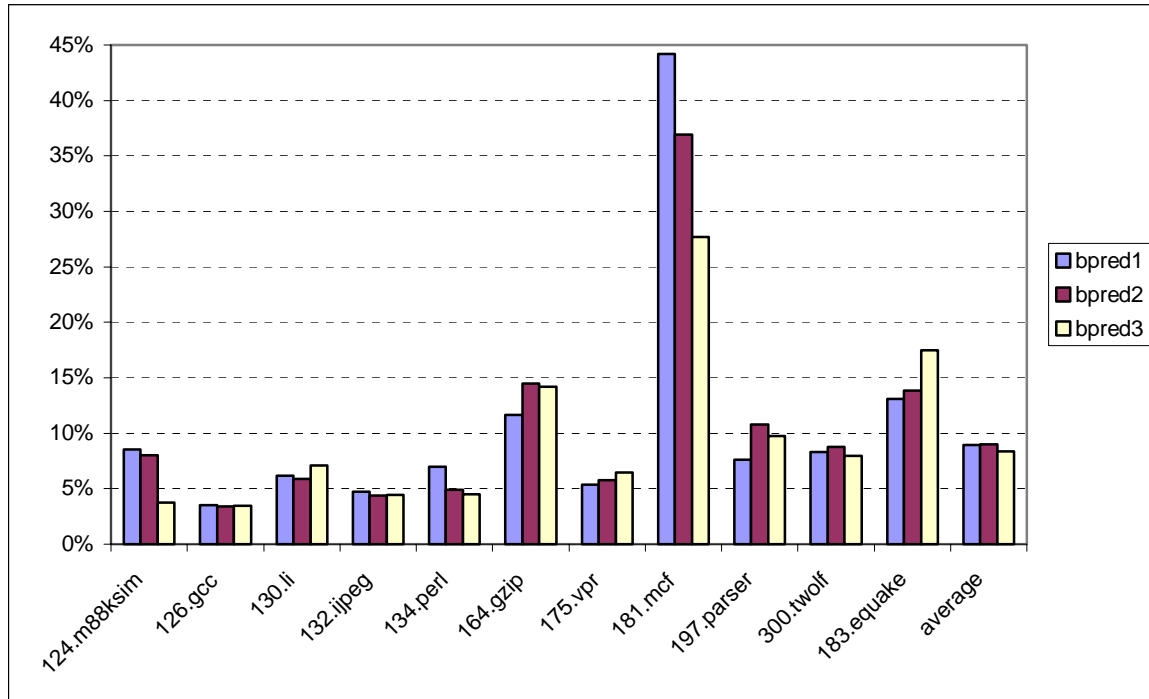


Figure 18. Speedup of the *wpc* compared to the *orig* configuration with various branch prediction accuracies.

5 RELATED WORK

There have been several studies examining how speculation affects multiple issue processors [1-8]. Farkas *et al* [1], for example, looked at the relative memory system performance improvement available from techniques such as non-blocking loads, hardware prefetching, and speculative execution, used both individually and in combination. The effect of deep speculative execution on cache performance has been studied by Pierce and Mudge [2]. Several other authors [3-8] have also examined speculation and pre-execution in their studies. A study by Wallace *et al.* [4] introduced instruction recycling, where previously executed wrong path instructions are injected back into the rename stage instead of being discarded. This increases the supply of instructions to the execution pipeline and decreases fetch latency.

Prefetching, which overlaps processor computations with data accesses, has been shown to be one of several effective approaches that can be used to tolerate large memory latencies. Prefetching can be hardware-based, software-directed, or a combination of both [26]. Software prefetching relies on the compiler to perform static program analysis and to selectively insert prefetch instructions into the executable code [19-24]. Hardware-based prefetching, on the other hand, requires no compiler support, but it does require some additional hardware connected to the cache [9-18]. This type of prefetching is designed to be transparent to the processor.

Jouppi [10] proposed *victim caching* to tolerate the conflict misses in the cache. While several other prefetching schemes have been proposed, such as adaptive sequential prefetching [11], prefetching with arbitrary strides [12-13, 17], fetch directed prefetching [16], and selective prefetching [18], Pierce and Mudge [25] have proposed a scheme called wrong path instruction prefetching. This mechanism combines next-line prefetching with the prefetching of

all instructions that are the targets of branch instructions regardless of the predicted direction of conditional branches.

Most of the previous prefetching schemes require a significant amount of hardware to implement. For instance, they require a *prefetcher* that prefetches the contents of the missed address into the data cache or into an on-chip prefetch buffer. Furthermore, a *prefetch scheduler* is needed to determine the right time to prefetch. On the other hand, this work has shown that executing loads down the wrongly-predicted branch paths can provide a form of indirect prefetching. However, this indirect prefetching may introduce some cache pollution. Our proposed Wrong Path Cache (WPC) is essentially a combination of a very small prefetch buffer and a victim cache [10] to eliminate this pollution effect.

6 FUTURE WORK

While this study has examined the effects of several parameters on the performance of the WPC, there are still many important factors left to be considered. For instance, increasing the processor's issue width will tend to cause more wrong-path loads to be executed. Furthermore, additional work still needs to be done to fully understand the relationship of the branch prediction accuracy to the performance of the WPC.

This study indicates that the execution of wrong-path loads has significant potential for improving memory access times in highly aggressive processor designs. The Alpha 21264 and the IBM Power4, for instance, both incorporate a load miss queue (LMQ) to keep track of L1 cache misses. Instructions are flushed from the pipeline when a branch misprediction occurs. However, entries in the LMQ are simply marked instead of being flushed. When the data returns for an entry that is marked in the LMQ, it is put into the L1 cache, but it is not written to the register file. In the WPC design presented in this paper, we go one step further than existing processors by intentionally continuing to execute as many ready loads as possible down the mispredicted branch path *even after the branch is resolved*, as long as there are memory ports available. A mechanism such as the LMQ can be used to determine whether a value should be put into the WPC after the miss is serviced rather than making this decision at issue time. This delay in determining what should be placed in the WPC should further increase its potential performance.

7 CONCLUSIONS

This study examined the performance effects of executing the load instructions that are issued along the incorrectly predicted path of a conditional branch instruction. While executing these wrong-path loads increases the total number of memory references, we find that allowing these loads to continue executing, even after the branch is resolved, can reduce the number of misses observed on the correct branch path. Executing these wrong-path loads thus provides an indirect prefetching effect. For small caches, however, this prefetching can pollute the cache causing an overall slowdown in performance.

We proposed the Wrong Path Cache (WPC), which is a combination of a small prefetch buffer and a victim cache, to eliminate the pollution caused by the execution of the wrong-path loads. Simulation results show that, when using an 8 KB L1 data cache, the execution of wrong-path loads without the WPC can result in a speedup of up to 5%. Adding a fully-associative eight-entry WPC to an 8 KB direct-mapped L1 data cache, though, allows the

execution of wrong path loads to produce speedups of 4% to 37% with an average speedup of 9%. The WPC also shows substantially higher speedups compared to the baseline processor equipped with a victim cache of the same size.

This study has shown that the execution of loads that are known to be from a mispredicted branch path has significant potential for improving the performance of aggressive processor designs. This effect is even more important as the disparity between the processor cycle time and the memory speed continues to increase. The Wrong Path Cache proposed in this paper is one possible structure for exploiting the potential benefits of executing wrong-path load instructions.

ACKNOWLEDGEMENTS

The authors would like to thank Peng-fei Chuang, Joshua J Yi, Chris Hescott, Baris M Kazar, and Keqiang Wu for their helpful comments on previous drafts of this paper.

This work was supported in by National Science Foundation grants EIA-9971666 and CCR-9900605, by IBM Corporation, by Compaq's Alpha Development Group, and by the Minnesota Supercomputer Institute. D. Lilja was supported by a Fulbright Award from the Australian-American Education Foundation during portions of this work.

REFERENCES

- [1] K. I. Farkas, N. P. Jouppi, and P. Chow, "How Useful Are Non-Blocking Loads, Stream Buffers, and Speculative Execution in Multiple Issue Processors?" *Technical Report WRL RR 94/8*, Western Research Laboratory – Compaq, Palo Alto, CA, August 1994.
- [2] J. Pierce and T. Mudge, "The effect of speculative execution on cache performance," *IPPS 94, Int. Parallel Processing Symp.*, Cancun Mexico, pp. 172-179, Apr. 1994.
- [3] G. Reinman, T. Austin, and B. Calder, "A Scalable Front-End Architecture for Fast Instruction Delivery," *26th International Symposium on Computer Architecture*, pages 234-245, May 1999.
- [4] S. Wallace, D. Tullsen, and B. Calder, "Instruction Recycling on a Multiple-Path Processor," *5th International Symposium On High Performance Computer Architecture*, pages 44-53, January 1999.
- [5] G. Reinman and B. Calder, "Predictive Techniques for Aggressive Load Speculation," *31st International Symposium on Microarchitecture*, pages 127-137, December 1998.
- [6] D. Lee, J.-L. Baer, B. Calder, and D. Grunwald, "Instruction Cache Fetch Policies for Speculative Execution," *22nd International Symposium on Computer Architecture*, pages 357-367, June 1995.
- [7] J. D. Collins, H. Wang, D. M. Tullsen, C. Hughes, Y.-F. Lee, D. Lavery, J. P. Shen, "Speculative Precomputation: Long-range Prefetching of Delinquent Loads," In *28th International Symposium on Computer Architecture*, July, 2001.
- [8] J. Dundas and T. Mudge, "Improving data cache performance by pre-executing instructions under a cache miss," *Proc. 1997 ACM Int. Conf. on Supercomputing*, July 1997, pp. 68-75.
- [9] A.J. Smith, "Cache Memories," *Computing Surveys*, Vol. 14, No. 3, Sept. 1982, pp. 473-530.
- [10] N.P. Jouppi, "Improving Direct-Mapped Cache Performance by the Addition of a Small Fully-associative Cache and Prefetch Buffers," *Proc. 17th Annual International Symposium on Computer Architecture*, Seattle, WA, May 1990, pp. 364-373.
- [11] F. Dahlgren, M. Dubois and P. Stenstrom, "Fixed and Adaptive Sequential Prefetching in Shared-memory Multiprocessors," *Proc. First IEEE Symposium on High Performance Computer Architecture*, Raleigh, NC, Jan. 1995, pp. 68-77.
- [12] J.W.C. Fu and J.H. Patel, "Data Prefetching in Multiprocessor Vector Cache Memories," *Proc. of the 18th International Symposium on Computer Architecture*, Toronto, Ont., Canada, May 1991, pp. 54-63.
- [13] T.F. Chen and J.L Baer, "Effective Hardware-Based Data Prefetching for High Performance Processors," *IEEE Transactions on Computers*, Vol. 44, No.5, May 1995, pp. 609-623.
- [14] D. Joseph and D. Grunwald, "Prefetching using markov predictors," *IEEE Transactions on Computers*, Vol. 48, No 2, 1999, pp 121-133.

- [15] A. Roth, A. Moshovos, and G.S. Sohi, "Dependence based prefetching for linked data structures," *In Proc. ASPLOS-VIII*, pp. 115–26, Oct. 1998.
- [16] G. Reinman, B. Calder, and T. Austin, "Fetch Directed Instruction Prefetching," *In proceedings of the 32nd International Symposium on Microarchitecture*, November 1999.
- [17] T.F. Chen and J.L. Baer, "A Performance Study of Software and Hardware Data Prefetching Schemes," *Proc. of the 21st Annual International Symposium on Computer Architecture*, Chicago, IL, April 1994, pp. 223-234.
- [18] R. Pendse and H. Katta, "Selective Prefetching: Prefetching when only required," *Proc. of the 42nd IEEE Midwest Symposium on Circuits and Systems*, volume 2, 2000, pp. 866-869.
- [19] T. C. Mowry, M. S. Lam, and A. Gupta, "Design and Evaluation of a Compiler Algorithm for Prefetching," *Proc. of the Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, October 1992.
- [20] C-K. Luk and T. C. Mowry. "Compiler-based prefetching for recursive data structures," *In Seventh International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 222--233, Oct. 1996.
- [21] Bernstein, D., C. Doron and A. Freund, "Compiler Techniques for Data Prefetching on the PowerPC," *Proc. International Conf. on Parallel Architectures and Compilation Techniques*, June 1995, pp. 19-26.
- [22] T. Mowry and A. Gupta, "Tolerating Latency through Software-controlled Prefetching in Shared-memory Multiprocessors," *Journal of Parallel and Distributed Computing*, Vol. 12, No. 2, June 1991, pp. 87-106.
- [23] E.H. Gornish, E.D. Granston and A.V. Veidenbaum, "Compiler-directed Data Prefetching in Multiprocessors with Memory Hierarchies," *Proc. 1990 International Conference on Supercomputing*, Amsterdam, Netherlands, June 1990, pp. 354-368.
- [24] M.H. Lipasti, W.J. Schmidt, S.R. Kunkel and R.R. Roediger, "SPAID: Software Prefetching in Pointer and Call-Intensive Environments," *Proc. 28th Annual International Symposium on Microarchitecture*, Ann Arbor, MI, November 1995, pp. 231-236.
- [25] J. Pierce and T. Mudge, "Wrong-Path Instruction Prefetching," *Proc. of 29th Annual IEEE/ACM Symp. Microarchitecture (MICRO-29)*, Dec. 1996, pp. 165-175.
- [26] S. P. VanderWiel and D. J. Lilja, "Data Prefetch Mechanisms," *ACM Computing Surveys*, Vol. 32, Issue 2, June 2000, pp. 174-199.
- [27] D.C. Burger, T.M. Austin, and S. Bennett, "Evaluating future Microprocessors: The SimpleScalar Tool Set," *Technical Report CS-TR-96-1308*, University of Wisconsin-Madison, July 1996.
- [28] C. Price, "MIPS IV Instruction Set, revision 3.1.," *MIPS Technologies, Inc.*, Mountain View, CA, January 1995.
- [29] AJ KleinOsowski, J. Flynn, N. Meares, and D. J. Lilja, "Adapting the SPEC 2000 Benchmark Suite for Simulation-Based Computer Architecture Research," *Workshop on Workload Characterization, International Conference on Computer Design*, Austin, TX, Sept 18-20, 2000.
- [30] S-T Pan, K. So, and J.T. Rahmeh, "Improving the Accuracy of Dynamic Branch Prediction Using Branch Correlation," *Proc. of the 5th International Conference on Architectural Support for Programming Languages and Operating Systems*, 1992, pp. 76-84.
- [31] D. A. Patterson and J. L. Hennessy: *Computer Architecture: A Quantitative Approach*, 2nd edition, Morgan Kaufmann press, 1995, pp. 393-395.
- [32] T.Y. Yeh and Y. N. Patt, "A Comparison of Dynamic Branch Predictors that Use Two Levels of Branch History," *Proc. of the International Symposium on Computer Architecture*, 1993, pp. 257--267.