

Feedback-Directed Data Cache Optimizations for the x86

Ronald D. Barnes* Ronnie Chaiken David M. Gillies

Microsoft Research

One Microsoft Way, Redmond WA, 98052

rdbarnes@crhc.uiuc.edu, {rchaiken, dgillies}@microsoft.com

Abstract

1 *The vast majority of desktop microprocessors in use*
2 *today belong to a single architectural family, the x86.*
3 *The success of this architecture has led to a large*
4 *number of microarchitectures and a growing need to*
5 *evolve the ISA to meet the changing demands of*
6 *applications. Unfortunately, most compiled code*
7 *today targets the 486 or Pentium[®],[‡] thereby missing*
8 *performance opportunities of newer processors. We*
9 *believe that specialized compilation within the x86*
10 *family can yield large performance gains over*
11 *generically compiled code. This paper examines the*
12 *effectiveness of the Pentium III data cache*
13 *management instructions on desktop applications. We*
14 *use a memory trace analysis in our optimization to*
15 *guide the placement of cache prefetch and cache*
16 *bypass instructions at the binary-level. Our results*
17 *show that significant performance improvements can*
18 *be achieved for a wide range of applications.*

1 Introduction

19 Since the advent of the 386, the number of
20 implementations within the x86 family of
21 microprocessors has grown very rapidly. These
22 implementations run the gamut of modern computer
23 architecture: non-pipelined to pipelined, scalar to
24 superscalar, in-order to out-of-order, cacheless to
25 having a highly stratified cache structure. At the same
26 time, the ISA has been evolving to meet the changing
27 demands of applications. ISA enhancements include
28 support for 32 bit addressing, partial instruction
29 predication, multimedia extensions with SIMD
30 support, and an increase and widening of the register
31 set. In addition, the circuitry surrounding the CPU has
32 changed. For example, the sizes of off-chip cache and
33 main memory have increased many-fold. Factoring
34 the contributions of several different manufacturers
35 into the mix has made for a very diverse landscape
36 within the x86 family.

37 Unfortunately, most code intended for the x86
38 platform is compiled targeting the 486 or Pentium
39 thereby missing performance opportunities of newer
40 processors. There are several reasons for this choice.

41 In order to deliver software that runs on the full range
42 of x86 platforms one cannot make assumptions about
43 the existence of special new instructions or try to
44 exploit a new microarchitectural feature. Also, due to
45 the ever-tightening product development cycle time,
46 testing multiple executables has generally prevented
47 machine-specific versioning. This limitation has been
48 overcome in a very small number of cases by the
49 conditional execution of specialized code within
50 dynamically linked libraries. These examples are
51 often handwritten in assembly language to maximize
52 the performance of a few critical routines. However,
53 this is far from a complete solution for general
54 programs.

55 While there are a number of pragmatic reasons for
56 adopting the one-size-fits-all approach to software
57 delivery, we contend that significant potential
58 performance gains offered by variations in x86
59 platforms are being overlooked. In this paper we
60 demonstrate one area for performance improvement in
61 this space. Specifically, we make use of the Pentium
62 III's data cache management instructions to illustrate
63 the gains that are available in general programs. Our
64 optimization heuristics improve data cache
65 performance by inserting prefetch and bypass
66 instructions directly into the binary, through binary
67 rewriting. The heuristics use information that is
68 gathered from a cache simulation that consumes
69 memory traces on the fly.

70 The remainder of this paper is divided into five
71 sections. Section 2 reviews previous related work in
72 this area. Next we describe the methodology used in
73 our experiments. Section 4 describes the Pentium III
74 instructions used in this study and the heuristics used
75 to drive the optimizations are described in Section 5.
76 Our results are presented in Section 6 and the final
77 section contains some concluding remarks.

2 Previous Work

78 Previous work advocating processor-specific
79 optimizations on x86 processors has been done by
80 Merten[13]. Much of this work focused on a
81 framework that enables optimizations, rather than on
82 the optimizations themselves. In addition, the
83 optimizations presented were mostly ad-hoc pattern
84 matching. In this work, we present a trace-directed
85 analysis to guide the placement of data cache
86 management instructions.

* Ronald D. Barnes is currently with the IMPACT Research Group in the Center for Reliable and High-Performance Computing, University of Illinois, Urbana-Champaign, Illinois 61801.

[‡] Intel, Pentium, MMX and Pentium[®] III Xeon are either registered trademarks or trademarks of Intel Corporation in the United States and/or other countries.

Table 1. Cache specifics of the target Pentium III Xeon machine.

| | |
|---------------------------------------|-------------------|
| L1 Cache Size | 16kB |
| L1 Associativity | 4-way |
| L1 access latency | 1 cycle (assumed) |
| L2 Cache Size (instruction & data) | 512kB |
| L2 access latency (measured) | ~20 cycles |
| Main memory latency (measured) | ~65 cycles |
| Cache Replacement | LRU (assumed) |
| Cache block size | 32 bytes |

1 The performance limitations caused by the ever-
 2 widening gap between processor and main memory
 3 speed are well understood, and hierarchical caches
 4 have been used to ameliorate this effect. Methods for
 5 instruction cache optimizations, like code
 6 reordering[12][21] and instruction cache
 7 prefetching[11], have been studied and shown to
 8 increase performance in some cases. Hardware
 9 mechanisms for runtime data cache management[8]
 10 have been examined; these mechanisms, however,
 11 often require expensive hardware support. Although
 12 dynamic scheduling has been shown to increase the
 13 tolerance of load latency[19], even machines with out-
 14 of-order execution benefit from prefetching[3] and
 15 load speculation[18].

16 Standard compiler techniques[15] for data
 17 prefetching are well established and are used in many
 18 modern compilers. These algorithms, however, tend
 19 to rely on strided loops such as those in scientific
 20 applications[17]. Since such loops are rare in integer
 21 programs, data prefetching is not often used for
 22 general applications. Much of the work[16] relating to
 23 prefetching in integer programs has relied upon
 24 informing memory operations that expose the cache
 25 behavior of instructions to the program. However,
 26 such instructions are not yet available in the x86
 27 microarchitectural family. In this work, we
 28 demonstrate that significant improvements can be
 29 achieved by using Pentium III specific cache
 30 management instructions and that this is applicable to
 31 a wide range of integer and floating point applications.

3 Trace Methodology

32 To effectively utilize the instruction set extensions, we
 33 rely on memory trace information to locate
 34 opportunities for optimizations. To gather this
 35 information, every memory reference in a binary is
 36 instrumented to generate a call to a runtime routine
 37 that generates output to a trace buffer with information
 38 about the memory reference. The instrumentation code
 39 sends the address of the memory instruction, the
 40 address of the accessed data and the corresponding
 41 size of the access to a runtime routine on each
 42 invocation. Control is occasionally passed to a cache
 43 simulator that consumes this trace on the fly. The
 44 simulator, which is modeled after the Pentium III
 45 Xeon™ L1 data cache, keeps track of the miss
 46 statistics for each memory reference. The specifics of

47 the cache hierarchy are shown in Table 1. The latency
 48 values are determined experimentally and are used
 49 later to determine the optimal prefetch distance.

50 To enable prefetch and bypass optimizations the
 51 simulation keeps track of memory accesses and their
 52 reuse patterns. Four types of information are gathered
 53 during the simulation that is used later in the
 54 optimization phase: 1. a record of the miss frequency
 55 of all loads and stores is kept; 2. the strides between
 56 successive dynamic occurrences of each memory
 57 instruction are recorded; 3. the result of analysis
 58 performed to determine if the cache lines brought in
 59 by each memory instruction tend to be reused before
 60 being replaced from the cache is stored; 4. lines that
 61 are displaced by store instructions are recorded to
 62 determine if those lines are later brought back into the
 63 cache on a miss. This last mechanism detects
 64 situations where the miss could be avoided if the
 65 earlier store instruction had not write-allocated.

66 During the analysis phase, the executables and
 67 dynamically linked libraries of the benchmarks are
 68 instrumented. The C runtime libraries however are
 69 not instrumented, nor are the Microsoft® Windows
 70 NT® system libraries. While the simulation thus sees
 71 only a partial list of all data memory accesses, the
 72 number of memory references and cache misses found
 73 during the simulation was a close approximation when
 74 measured against the Pentium III performance
 75 counters. This validation check provides some level of
 76 confidence in our simulation.

77 Both the instrumentation and optimization of the
 78 benchmarks is performed at the binary level using a
 79 post-link binary rewriting technology called Vulcan
 80 [20]. This technology is similar to ATOM[4],
 81 EEL[10], or the IMPACT Binary Reoptimization
 82 System[13]. Information from the simulation is fed to
 83 a post-link optimizer tool to produce an optimized
 84 binary using the heuristics detailed in Section 5. The
 85 optimized binaries are then rerun using the same input
 86 data as the training set to measure the effectiveness of
 87 the transformation. All of our measurements are made
 88 using a Pentium III Xeon 500 MHz machine running
 89 Microsoft Windows NT 4.0. Each program is run
 90 repeatedly on an unloaded machine and the actual
 91 execution time is measured. The additional

* Microsoft and Windows NT are registered trademarks of Microsoft Corporation in the United States and/or other countries.

Table 2 – Data cache control instructions

| | |
|-------------|-----------------------------------------------------------------------------------|
| prefetcht0 | Prefetch into 0 th level cache (Into both L1 and L2 in Pentium III) |
| prefetcht1 | Prefetch into 1 st level cache (Into just L2 in Pentium III) |
| prefetcht2 | Prefetch into 2 nd level cache (Into just L2 in Pentium III) |
| prefetchnta | Prefetch non-temporal data (Into just L1 in Pentium III) |

prefetch instructions

| | |
|---------|--------------------------------------------------------------------------|
| movntps | 128 bit streaming store (From Streaming SIMD Extensions register set) |
| movntq | 64 bit streaming store (From MMX register set) |
| maskmov | Masked variable length streaming store (From MMX register set) |

streaming store instructions

1 performance metrics detailed in Section 6 are also
 2 measured during these runs using the Pentium III
 3 performance counters.
 4 Although all the described optimizations are
 5 performed on existing binaries, the techniques
 6 described can be used at compile time to generate a
 7 machine-specific executable. In addition, since the
 8 source code is not used, these optimizations possibly
 9 could be done on the host machine. This might be
 10 facilitated by hardware profiling mechanisms[4][14].

4 Using the Pentium III Instructions

11 The Intel® Pentium III microarchitecture features
 12 several ISA extensions that can be used to manage
 13 placement of data in the memory hierarchy. The
 14 extensions include both data cache prefetching and
 15 cache-bypassing types of instructions. These
 16 instructions can dramatically improve the performance
 17 of the memory hierarchy, and thus can also
 18 substantially improve application performance.

4.1 ISA Extensions

19 The Pentium III's cache control instructions are shown
 20 in Table 2. Four prefetch instructions are provided to
 21 allow prefetching into different levels of the cache
 22 hierarchy. Since only two levels of cache are visible
 23 on the Pentium III, the effect of the `prefetcht1`
 24 and `prefetcht2` instructions is the same. Three
 25 streaming store instructions provide support for stores
 26 of various sizes that bypass the cache hierarchy. Each
 27 of these instructions performs a non-allocating store
 28 from either the 64 or 128 bit registers. These stores
 29 bypass the cache hierarchy if their cache line is not
 30 already present in the cache, otherwise they act as
 31 normal store instructions. The `maskmov` instruction
 32 performs a streaming store of selected bytes from an
 33 MMX register that are determined by a mask in an
 34 additional MMX register. For more information about
 35 these instructions see [6] and [7]. Further information
 36 on how these instructions are used in this work, and
 37 the difficulties and costs for employing them can be
 38 found in [1].

4.2 Data Prefetching

39 A significant body of previous research demonstrates
 40 that software controlled data prefetching can

41 significantly improve cache effectiveness and system
 42 performance (see Section 2). In order to effectively
 43 use any type of data prefetching, however, it is
 44 necessary to determine how far in advance of a load to
 45 place a corresponding prefetch. To determine this
 46 optimal distance, experimental measurements are
 47 made. The results of using prefetch instructions in this
 48 test program executing on a Pentium III are presented
 49 in Figure 1. In at test program, the number of cycles
 50 between the prefetch and a load are varied from 1 to
 51 100. The program consisted of long sequences of
 52 dependent instructions to minimize the effect of the
 53 out-of-order execution of the Pentium III. The test
 54 program also ensures that the load targets would reside
 55 in the L2 cache, but would not reside in the L1 cache.
 56 Each run of the test program is identical, except for
 57 the number of dependant instructions between the load
 58 and store. Figure 1 shows that the optimal number of
 59 cycles needed for prefetching an L1 cache miss in our
 60 test program is around 18 cycles. Not surprisingly,
 61 this corresponds approximately to the latency of a L1
 62 cache miss as given in Section 3. It is important to
 63 note from Figure 1 that as long as the prefetch is
 64 issued a few cycles before the corresponding load,
 65 some performance gain is achieved, and this benefit
 66 increases relatively linearly up to the optimal distance.
 67 However, if a prefetch occurs too far in advance of the
 68 load, there is an increased chance that the desired
 69 cache line might be displaced before it is needed or the
 70 fetched line might displace some data that is needed
 71 first. When the load must go to main memory,
 72 maximum benefit will come from prefetching much
 73 farther in advance. Thus, maximum benefit comes
 74 from knowing which level of the memory hierarchy
 75 the data resides and scheduling the prefetch
 76 appropriately.

4.3 Important issues in using cache control instructions

77 The use of the machine-specific optimizations
 78 considered in this work is not without cost. In the case
 79 of data prefetching, the inserted prefetches can
 80 significantly increase the dynamic instruction count.
 81 If the memory hierarchy is swamped with prefetches it
 82 can cause trailing loads to stall. Also, extraneous
 83 prefetches can displace useful cache lines. When
 84 streaming stores are performed, an increase in memory
 85 transactions can occur if the cache line of the stored
 86 value is loaded soon after the store. In addition, since

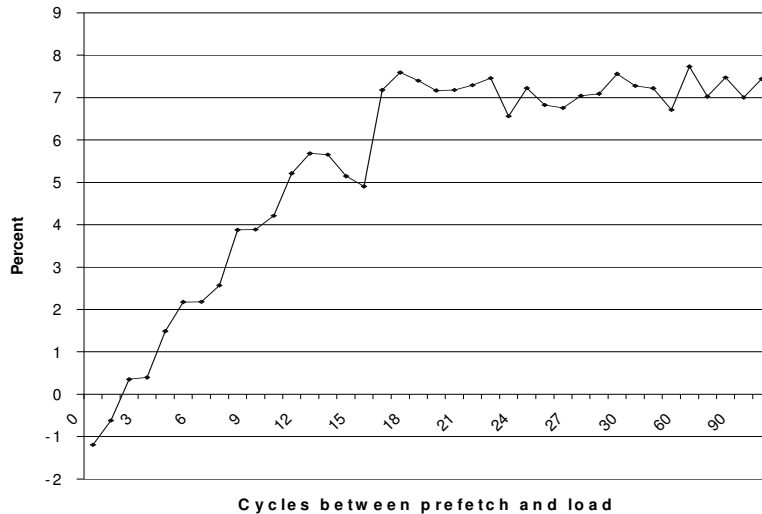


Figure 1. Speedup vs. prefetch distance in an experimental program.

1 there are no streaming stores from the 32-bit register
 2 set, general purpose stores cannot be made to be
 3 streaming without performing potentially costly
 4 moves of the data to the MMX™ register set. In
 5 addition, to do a 32-bit store requires setting up an
 6 additional mask register for the `maskmov` instruction.

5 Code Transformations and Heuristics

7 To effectively use the ISA extensions during
 8 optimization, several heuristics are used for
 9 determining where and when to apply transformations
 10 that use these instructions. The heuristics are tuned to
 11 minimize the total number of transformations made to
 12 the code by optimizing only those memory accesses
 13 that have a high probability of causing a cache miss.
 14 This might leave some opportunities unexplored but
 15 the goal of this work is to demonstrate the feasibility
 16 of such machine specific optimizations. This section
 17 details some of the code transformations and
 18 optimization heuristics that are used in this work. For
 19 further details see [1].

5.1 Important issues in using cache control instructions

20 There are several kinds of code transformations that
 21 are useful for making the optimizations considered in
 22 this work more effective and less expensive. This
 23 usually requires transforming the code sequence to
 24 make it more amenable to the optimization.

25 Often randomly strided memory accesses occur
 26 inside loops for integer applications. In these cases, it
 27 is often not possible to calculate the target address of
 28 loads in future iterations. One technique that might
 29 help in this situation is to speculatively prefetch across
 30 the back edge of a loop. In such cases, loop
 31 rotation[15] may allow more distance to be placed
 32 between the prefetch and its matching load. As shown

33 in Figure 2, this technique allows the target address
 34 calculation for a load in the next iteration to be
 35 speculatively performed in the current iteration. In
 36 comparison to the original code sequence, this type of
 37 speculation causes more instructions to be executed, as
 38 it performs the address calculation for one more
 39 iteration than the original code. Therefore, it is
 40 important to ensure that the loop iteration count is
 41 relatively large so that the additional calculation
 42 becomes relatively insignificant.

43 Frequently, the address calculation for a candidate
 44 load instruction is dependent on another load
 45 instruction. In this case, the control flow often limits
 46 the distance that the two loads can be moved apart;
 47 thereby limiting the distance between the prefetch and
 48 the load. However, using conditional move
 49 instructions, the load can be moved above a branch
 50 without breaking the semantics of the program, as the
 51 execution of the load is still predicated on the branch
 52 condition. Now the corresponding prefetch can move
 53 above the branch, and the number of cycles between
 54 the issue of the prefetch and the load can be increased,
 55 improving the effectiveness of the optimization. This
 56 type of speculative load transformation is shown in
 57 Figure 3.

58 The x86 architecture features string instructions that
 59 can load and store to memory using a single
 60 instruction. Often these instructions are designated
 61 with a `rep` prefix so that they are repeated without
 62 fetching and decoding additional instructions. In this
 63 way, a single instruction can operate on a large
 64 amount of memory. These instructions are found in
 65 many general programs and are amenable to special
 66 optimizations. Since these instructions access a large
 67 portion of contiguous memory, the 64-bit registers can
 68 be used to operate on the majority of this data,
 69 reducing the dynamic number of memory operations.
 70 The single instruction is replaced with an optimized
 71 loop that can utilize both prefetches and streaming
 72 stores from the MMX register set. The string
 73 instructions use a loop counter to specify how much

1 data is to be accessed. This counter can be properly
2 adjusted and the majority of the data can be accessed
3 through 64-bit MMX registers, facilitating the use of
4 streaming stores of this data. For string instructions
5 with poor cache behavior, performance is greatly
6 improved by using prefetches to bring data needed for
7 future iterations of the loop into the cache and by
8 using streaming stores to store data with poor temporal
9 locality. Although these optimizations increase the
10 number of instructions that execute, they improve
11 execution performance.

5.2 Prefetch Heuristics

12 Much of the prior work related to data prefetching has
13 looked at optimizing array accesses in inner loops of
14 numerical applications. Unfortunately, these accesses
15 occur infrequently in general integer applications. To
16 determine if a load behaves in a strided fashion, our
17 simulation uses a form of value profiling[2][9] in the
18 simulation on the target addresses of the load
19 instructions. For every static load instruction in an
20 executable the simulation maintains a last address-
21 accessed datum. The difference between the current
22 memory location and the last memory location is a
23 stride that gets recorded into a stride list for a
24 particular load instruction. For a load that accounts
25 for more than five percent of all the misses, we use the
26 most commonly encountered strides to determine the
27 prefetch optimization.

28 When a candidate has a single dominant stride
29 (>70%) we found that placing it in the cache as early
30 as possible gives the best results. Intel defines the
31 number of iterations to prefetch ahead as the *prefetch*
32 *distance* and provides a formula for its calculation[7].
33 A similar, but slightly simpler approach is taken in this
34 work. The length of a loop iteration in dynamic
35 instructions is used to determine the distance to
36 prefetch ahead. Since the loop may contain some
37 control flow, profiling information can be used to
38 determine the frequency of each control flow edge,
39 and the average length of a loop iteration is used. The
40 prefetch is inserted to attempt to prefetch 25 dynamic
41 instructions before the load. Thus the number of
42 iterations to prefetch ahead is equal to 25 divided by
43 the average length of a loop in instructions. The
44 distance of 25 dynamic instructions was chosen
45 experimentally, and corresponds to an IPC of
46 approximately 1.4 if the distance in cycles is desired to
47 be 18.

48 Often in integer programs, memory accesses fail to
49 follow any strided pattern. Only the current iteration
50 of a loop is prefetched in these cases. It is generally
51 difficult to prefetch effectively within one loop
52 iteration. However, several of the optimizations
53 described earlier in this section can be used to increase
54 the distance between the availability of the load
55 address and the execution of the load instruction. A
56 candidate falling into this category is not prefetched if
57 the prefetch cannot be placed at least three instructions
58 before the load, since as shown in Figure 1, inserting

59 prefetches immediately before corresponding loads
60 hurts performance.

5.3 Streaming Store Heuristics

61 Using the streaming store instructions in general
62 programs is slightly more complicated than using the
63 prefetching instructions. There is the additional cost
64 described in of moving data to the MMX register set
65 before performing the store that must be accounted for
66 in the cost analysis.

67 Candidates for streaming store optimization are
68 chosen via the same cache simulation used for the
69 prefetch analysis. We found that the best candidates
70 for optimization are those stores that miss in the cache
71 more than 90% of the time. Due to the cost of this
72 optimization, it is detrimental to use a streaming store
73 on a location that is already resident in the cache.
74 Cache hits due to cache lines brought in by previous
75 executions of the same instruction are not considered
76 hits in this determination. This helps to eliminate the
77 cache hits that would not occur if it a non-allocating
78 write had been executed. Of these candidates, only
79 the store instructions with data that are written back to
80 L2 before being read in more than 90% of its dynamic
81 occurrences are considered. Analysis is done to
82 determine if bringing a cache line in for each store
83 miss pollutes the cache and causes additional misses.
84 The number of additional misses caused by an
85 instruction is multiplied by the latency of a L1 cache
86 miss. This value is compared to the number of times
87 the instruction executed times 6 to account for the
88 overhead of inserting the code transformation. If the
89 effect of the additional cache misses is larger than the
90 assumed transformation cost, the benefit of using the
91 streaming store instruction likely outweighs the large
92 cost in the general case.

93 The cost of using the streaming store optimization in
94 general programs makes it unlikely to be of much
95 benefit. However, several special cases do occur in
96 general programs that make effective use of the
97 streaming store optimization. For example, if two
98 instructions store to consecutive memory locations and
99 these operations are both cache polluting then the two
100 stores can be combined into a single larger store using
101 the 8-byte registers. This has the advantage of possibly
102 reducing bus traffic and reducing the number of
103 instructions executed.

104 Another example of stores that can often be
105 effectively transformed is the string instructions. As
106 described earlier in this section, using the streaming
107 stores for this case is inexpensive since the code
108 transformation has already placed the data in an MMX
109 register. If the string instruction is likely to cause a
110 large number of cache misses and poor temporal
111 locality is identified, then expanding the instruction
112 into a small loop using both prefetching and streaming
113 stores tends to improve the performance.

```

LOOP:  MOV     EDX, DWORD PTR [ECX]
      ADD     EDX, ESI
      XOR     ECX, EDX
      ADD     EDI, DWORD PTR [ECX]
      .
      .
      .
      CMP    DWORD PTR [ECX], EBX
      JNE    LOOP

      MOV     EDX, DWORD PTR [ECX]
      ADD     EDX, ESI
      XOR     EAX, EDX
LOOP:  MOV     ECX, EAX
      MOV     EDX, DWORD PTR [ECX]
      ADD     EDI, EDX
      ADD     EDX, ESI
      XOR     EAX, EDX
      PREFETCHT0 BYTE PTR [EAX]
      .
      .
      .
      CMP    DWORD PTR [ECX], EBX
      JNE    LOOP

```

Figure 2. Using loop rotation to facilitate prefetching across a back edge

```

      .
      .
      .
      CMP    EAX, 0
      JE     NULL
      .
      .
      .
      MOV    EBX, DWORD PTR [EAX]
      .
      .
      .
      MOV    EDX, DWORD PTR [EBX]

      MOV    EBX, 0
      CMP    EAX, 0
      CMOVNE EBX, DWORD PTR [EAX]
      PREFETCHT0 BYTE PTR [EBX]
      .
      .
      .
      JE     NULL
      .
      .
      .
      MOV    EDX, DWORD PTR [EBX]

```

Figure 3. Using speculative load to increase distance between prefetch and load

6 Experimental Evaluation

1 By improving the cache behavior of application
 2 programs, significant overall performance increases
 3 can be achieved. To demonstrate that machine
 4 specific optimizations on the Pentium III can yield this
 5 improvement, seven test cases are run on six different
 6 applications. Each is probed and simulated to
 7 determine locations for optimization, and then
 8 optimized using a binary rewriting tool.

6.1 Benchmarks

9 The benchmarks in the study are chosen to represent a
 10 wide range of applications. *Compress* and *tomcatv* are
 11 taken from the SPEC95 integer and floating-point
 12 benchmark suites respectively. For these benchmarks,
 13 the SPEC reference inputs are used for performance
 14 measurements. *Ghostscript* is taken from the Aladdin
 15 Ghostscript 5.5 public release from the University of
 16 Wisconsin. *Ghostscript* is a postscript interpreter, and
 17 as the test case, a postscript file containing a large
 18 document is rastered and displayed. *Microsoft FoxPro*
 19 is a database application, and a scenario performing
 20 numerous transactions on a database is used as its test
 21 case. *Microsoft Word* and *Microsoft Excel* are large
 22 desktop applications, taken from the Microsoft Office
 23 2000 suite. As a test case for *word*, a large word
 24 document is run through the find-and-replace
 25 operation. For *excel*, two different test cases are used.
 26 The first test case stresses the recalculation engine of
 27 *excel* by making numerous changes to a large

28 spreadsheet. The second test case stresses the column
 29 editing operation over the same spreadsheet.

6.2 Performance

30 As shown in Figure 4, significant gain can be achieved
 31 on a wide range of applications. An average speedup
 32 of 6.8% with speedups as high as 27% is measured
 33 using the prefetch and bypass optimizations. Ignoring
 34 the high and low outlier results, an average of 4.3%
 35 speedup is obtained.

36 While the speedups for each benchmark vary
 37 considerably, only *compress* shows insignificant
 38 performance improvement. The performance statistics
 39 in Figure 4 is broken into three components: the
 40 combined effect, the prefetching optimization alone
 41 and the effect of the streaming store optimization
 42 alone*. Where the streaming store optimization is
 43 applied alone, a small decrease in performance is
 44 observed in all cases except *ghostscript*. There are a
 45 number of reasons for this decrease in performance.
 46 For *compress*, the insertion of several consecutive
 47 streaming store instructions has a serializing effect on
 48 instruction issue, since each of these instructions
 49 require the single complex instruction decoder
 50 available on the Pentium III. For the other
 51 applications, the weight of the enabling instructions
 52 for the streaming store optimization proves costly.
 53 The prefetching optimization alone has a positive

* *foxpro* and *tomcatv*, however, did not have streaming store optimization candidates.

1 effect on all applications. For *word* we have the
 2 situation where the combined optimizations yield a
 3 superior result to prefetching alone, even though the
 4 streaming store optimization alone has a negative
 5 effect. In the combined optimization, the use of
 6 prefetches, by placing greater pressure on the L1
 7 cache, has made the use of streaming stores more
 8 effective.
 9 *Word* is a clear outlier with the most improved
 10 performance. The majority of this improvement
 11 comes from the string optimization described in
 12 Section 5. Since a large number of changes are made
 13 to the *word* document in the test case, a routine for
 14 moving memory is called frequently. With the
 15 improved performance of this routine, the
 16 performance of *word* improves dramatically.
 17 Another way of evaluating the effectiveness of the
 18 optimizations is the utilization of the data cache as
 19 seen from Figure 5. This chart shows the weighted
 20 number of cycles that the processor might have to wait
 21 on the data cache to service an outstanding miss due to
 22 a load. This is calculated in the following way: during
 23 each cycle a count of the number of outstanding load
 24 misses is added to a running sum which is then
 25 normalized by the total number of cycles to execute
 26 the original program. For *tomcatv*, *compress*, *foxpro*
 27 and the *excel* insertion scenario, this number is
 28 originally very close to one. Thus, on average there is
 29 one outstanding cache miss throughout the run of the
 30 test. This is further evidence that the out-of-order
 31 execution of the Pentium III is quite successful in

32 tolerating load misses, since for this average to be
 33 close to one, often the processor must have several
 34 outstanding loads at a time. However, when
 35 optimized, there is a dramatic difference in this
 36 average. This difference is caused by the data
 37 prefetches that are started earlier than the load. These
 38 measurements are made using the performance
 39 monitoring counters of the Pentium III. Unfortunately
 40 this metric does not count any cycles for a load whose
 41 address is prefetched, even if the prefetching is not
 42 done in time. Thus the decrease in relative weighted
 43 number of cycles is the load latency that prefetching
 44 attempts to tolerate, not the latency that is actually
 45 tolerated.

6.3 Memory Bandwidth

46 An undesirable side effect of data prefetching is the
 47 increase in bus traffic as shown in Figure 6. This can
 48 happen if the prefetch displaces data from the cache
 49 that is used before the prefetched data, or if
 50 speculative prefetches are made for memory locations
 51 that are not in fact loaded. One notable result is the
 52 relatively dramatic increase in the memory bus
 53 transactions in the *ghostscript* benchmark. As the
 54 number of bus transactions almost double, this likely
 55 hampers some of the gain from using the cache
 56 optimization. It should also be noted however that
 57 while there is a large percentage increase in this bus
 58 traffic, the total number of memory bus transactions in
 59 *ghostscript* is actually rather small.

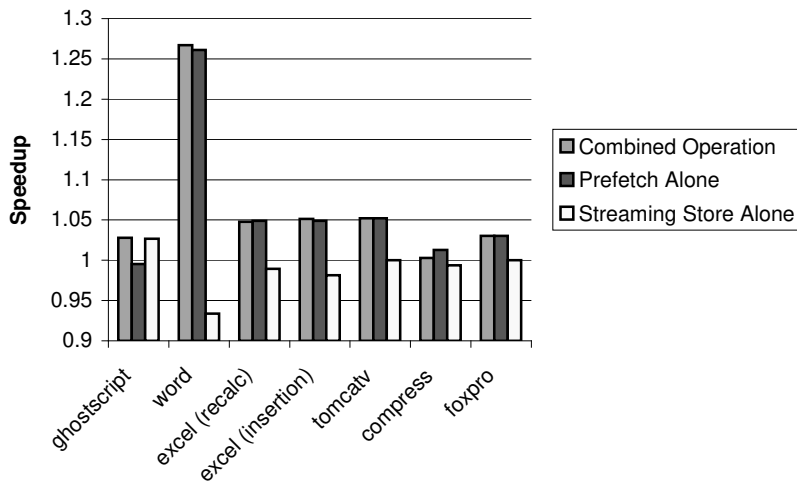


Figure 4. Breakdown of performance improvement

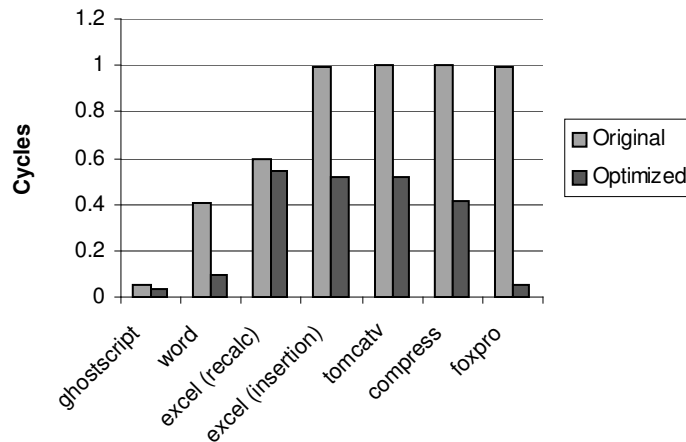


Figure 5. Weighted numbers of cycles with an outstanding cache miss

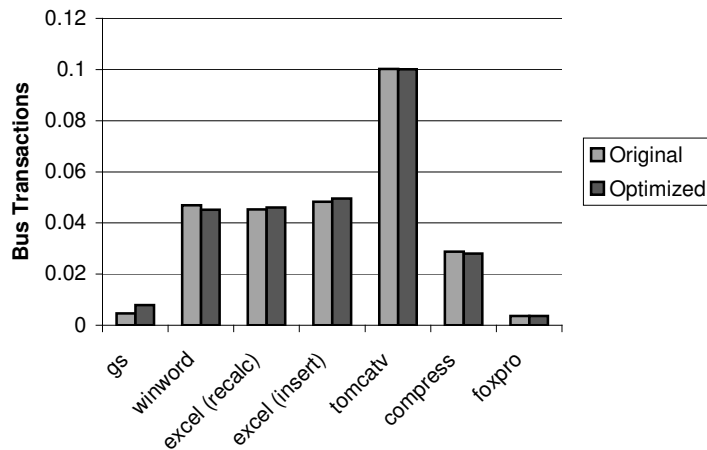


Figure 6. Bus transactions normalized by the original number of memory references

1 Since a streaming store is not write-allocating, it
 2 prevents the need for an initial read to bring the line
 3 into the cache hierarchy. This can serve to decrease
 4 the number of memory bus requests, canceling out
 5 some of the effect of increased pressure caused by
 6 data prefetching. In addition, the reduction in cache
 7 pollution and corresponding reduction in cache misses
 8 can reduce the traffic to memory as well. This
 9 positive effect is shown most prominently in Word. In
 10 addition, optimizations such as the string
 11 optimizations detailed in Section 4, use a 64-bit
 12 register where 32-bit registers were used previously.
 13 Since the memory bus of the Pentium III is 64 bits
 14 wide, this can also serve to reduce the pressure on the
 15 memory bus.

7 Conclusions

16 This paper provides some experimental support for
 17 specialized compilation within the x86 processor
 18 family. To do so, this paper studies the effectiveness

19 of the Pentium III prefetch and streaming store data
 20 cache-controlling instructions. These instructions
 21 are highly effective in optimizing the cache
 22 behavior under ideal circumstances. To show that
 23 the instructions can be useful to general programs,
 24 several applications are studied including non-loop
 25 intensive, integer desktop applications. Memory
 26 traces are collected with the aid of program
 27 articulation at the binary level and are concurrently
 28 processed by a cache simulator. With information
 29 from the cache simulation, heuristic optimization
 30 techniques are applied to identify and optimize,
 31 again at the binary level, a small number of
 32 instructions with poor cache behavior. The
 33 performance results show that the Pentium III
 34 cache-controlling instructions can be effectively
 35 utilized to achieve performance improvements in
 36 the range of 3-27%, with an average of 6.8%. With
 37 this and other specialized compilation techniques,
 38 the performance of applications within the x86

1 family can be greatly improved over generically
2 compiled code.

3 This paper also demonstrates that, for general
4 programs, it is relatively easy to utilize the Pentium III
5 prefetch instructions but that there can be significant
6 overhead associated with the use of the streaming
7 store instructions. To enable streaming store
8 optimization on a wide range of programs it is highly
9 recommended that the hardware implement a version
10 of the streaming store instruction using a general-
11 purpose register.

Acknowledgements

12 The authors would like to thank all the members of the
13 Programmer Productivity Research Center at
14 Microsoft Research and the IMPACT compiler team
15 at the University of Illinois for their support,
16 comments, and suggestions. We would also like to
17 thank the anonymous referees for their constructive
18 comments.

References

- 19 [1] Barnes, Ronald, Ronnie Chaiken, and David Gillies.
20 "Feedback-Directed Data Cache Optimizations for the
21 x86", Technical Report, Microsoft Research, Sept.
22 1999.
23 [2] Calder, Brad and Peter Feller, and Alan Eustace.
24 "Value Profiling and Optimization", *Journal of*
25 *Instruction Level Parallelism*, March 1999.
26 [3] Chen, William Y., Scott A. Mahlke, Pohua P. Chang,
27 and Wen-mei W. Hwu. "Data Access
28 Microarchitectures For Superscalar Processors with
29 Compiler-Assisted Data Prefetching", *Proceedings of*
30 *the 24th International Symposium on*
31 *Microarchitecture*, Nov. 1991.
32 [4] Dean, J., J. E. Hicks, C. A. Waldspurgen, W. E. Weihl,
33 and G. Chrysos, "Profileme: Hardware support for
34 instruction level profile-driven compilation using the
35 profile buffer," in *Proceedings of the 29th International*
36 *Symposium on Microarchitecture*, Dec. 1996.
37 [5] Eustace, Alan and Amitabh Srivastava. "ATOM: A
38 Flexible Interface for Building High Performance
39 Program Analysis Tools", USENIX Winter
40 Conference, 1995. Also available as WRL Technical
41 Report TN-44.
42 [6] *Intel Architecture Software Developer's Manual,*
43 *Volume 2: Instruction Set Reference.* Intel
44 Corporation, 1999.
45 [7] *Intel Architecture Software Optimization Reference*
46 *Manual.* Intel Corporation, 1999.
47 [8] Johnson, Teresa L., Matthew C. Merten, and Wen-mei
48 W. Hwu. "Run-time Spatial Locality Detection and
49 Optimization", *Proceedings of the 30th International*
50 *Symposium on Microarchitecture*, Dec. 1997.
51 [9] Kalamatianos, John, Ronnie Chaiken, and David
52 Kaeli. "Parameter Value Locality of Windows NT-
53 based applications," Workshop on PC-Performance
54 and Analysis held in conjunction with Micro-31,
55 November 1998.
56 [10] Larus, James R. and Eric Schnarr. "EEL: Machine-
57 Independent Executable Editing", *Proceedings of the*
58 *Conference on Programming Language Design and*
59 *Implementation*, June 1995.

- 60 [11] Luk, Chi-Keung and Todd C. Mowry.
61 "Cooperative Prefetching: Compiler and Hardware
62 Support for Effective Instruction Prefetching in
63 Modern Processors", *Proceedings of the 31st*
64 *International Symposium on Microarchitecture*,
65 Nov. 1998.
66 [12] McFarling, Scott. "Program Optimization for
67 Instruction Caches," *Proceedings of the 3rd*
68 *International Conference on Architectural Support*
69 *for Programming Languages and Operating*
70 *Systems*, April 1989.
71 [13] Merten, Mathew C. "A Framework for Profile-
72 Driven Optimization in the IMPACT Binary
73 Reoptimization System", Masters Thesis,
74 University of Illinois, 1999.
75 [14] Merten, Matthen C., Andrew R. Trick, Christopher
76 N. George, John C. Gyllenhaal, and Wen-mei W.
77 Hwu. "A Hardware-Driven Profiling Scheme for
78 Identifying Program Hot Spots to Support Runtime
79 Optimization", *Proceedings of the 26th*
80 *International Symposium on Computer*
81 *Architecture*, May, 1999.
82 [15] Muchnick, Steven S. *Advanced Compiler Design*
83 *& Implementation*, Morgan Kaufmann Publishers,
84 Inc., San Francisco, CA 1997.
85 [16] Mowry, Todd C. and Chi-Keung Luk. "Predicting
86 Data Cache Misses in Non-Numeric Applications
87 Through Correlation Profiling", *Proceedings of the*
88 *30th International Symposium on*
89 *Microarchitecture*, Dec. 1997.
90 [17] Mowry, Todd C., Monica S. Lam, and Anoop
91 Gupta. "Design and Evaluation of a Compiler
92 Algorithm for Prefetching", *Proceedings of the 5th*
93 *International Conference on Architectural Support*
94 *for Programming Languages and Operating*
95 *Systems*, October 1992.
96 [18] Reinman, Glenn and Brad Calder. "Predictive
97 Techniques for Aggressive Load Speculation",
98 *Proceedings of the 31st International Symposium on*
99 *Microarchitecture*, Nov. 1998.
100 [19] Srinivasan, Srikanth T., and Alvin R. Lebeck.
101 "Load Latency Tolerance In Dynamically
102 Scheduled Processors", *Proceedings of the 31st*
103 *International Symposium on Microarchitecture*,
104 Nov. 1998.
105 [20] Srivastava, Amitabh, et. al. "Vulcan", Technical
106 Report, Microsoft Research TR-99-76, Sept. 1999.
107 [21] Torrellas, J. and C. Xia, and R. Daigle. "Optimizing
108 Instruction Cache Performance for Operating
109 System Intensive Workloads", *IEEE Transactions*
110 *on Computers*, IEEE, Inc., Vol. 47, Number 12,
111 December 1998.