

How Useful Are Non-blocking Loads, Stream Buffers and Speculative Execution in Multiple Issue Processors?

Keith I. Farkas[†]
farkas@eecg.toronto.edu
416-978-1653

Norman P. Jouppi[‡]
jouppi@pa.dec.com
415-617-3305

Paul Chow[†]
pc@eecg.toronto.edu
416-978-2402

[†]Dept. of Electrical and Computer
Engineering
University of Toronto
10 Kings College Road
Toronto, Ontario, Canada
M5S 1A4

[‡]Digital Equipment Corporation
Western Research Lab
250 University Avenue
Palo Alto, California 94301

Abstract

We investigate the relative performance impact of non-blocking loads, stream buffers, and speculative execution both used individually and in conjunction with each other. We have simulated the SPEC92 benchmarks on a statically scheduled quad-issue processor model, running code from the Multiflow compiler. Non-blocking loads and stream buffers both provide a significant performance advantage, and their combination performs significantly better than either alone. For example, with a 64-byte, 2-way set associative cache with 32 cycle fetch latency, non-blocking loads reduce the run-time by 21% while stream-buffers reduce it by 26%, and the combined use of the two yields a 47% reduction. The addition of speculative execution further improves the performance of the systems that we have simulated, with or without non-blocking loads and stream buffers, by an additional 20% to 40%. We expect that the use of all three of these techniques will be important in future generations of microprocessors.

1 Introduction

A continuing trend in the design of computer systems is the widening gap between microprocessor and memory speeds. This speed discrepancy can have a significant impact on the performance of a system since it increases the time cost of servicing data cache misses. The performance impact can be lessened through the use of techniques that reduce the amount of time the processor is stalled for cache misses. In this paper, we consider non-blocking loads,

stream buffers, and speculative loads. We investigate the relative merits of these three techniques when used individually and in conjunction with each other, in the context of an advanced quad-issue statically scheduled microprocessor.

Non-blocking loads reduce the time stalled due to cache misses by allowing the processor to overlap the servicing of a miss with the execution of other instructions. The amount of overlap depends on the number of instructions that are available for execution that do not use the register being targeted by the load instruction. When an instruction is encountered that depends on the value being loaded and the load is still in-progress, the processor must stall until the load completes; a *true-data dependency* is said to exist between these two instructions. With a lockup cache, a stall will also occur if during the processing of a cache miss, other load or store instructions are executed. Such stalls can be avoided by using a *lockup-free cache*. A lockup-free cache allows multiple concurrent cache hits, misses, or both. By overlapping the processing of cache misses, the average time to service a miss decreases because the processor will be stalled for less time and more loads will complete in this time.

Prefetching also can reduce the time stalled due to cache misses. The goal of prefetching is to bring data closer to the processor before the processor requires it, rather than, as is the case with non-blocking loads, waiting for a cache miss to occur before initiating a fetch for missing data. The prefetch is initiated by some triggering event. With software prefetching, the trigger is prefetch instructions that are inserted into the code by a sophisticated compiler or user [3, 11]. With hardware prefetching, hardware is used to determine when a prefetch might be useful. Software

prefetching has been most successful on numeric codes, while hardware prefetching can be used with all types of applications (including the operating system).

Examples of hardware prefetch techniques include Chen and Baer's lookahead PC reference prediction method [4] and stream buffers [8, 12]. We study stream buffers as we believe them to be simpler and less-invasive than the lookahead scheme. The lookahead scheme is complicated by the need for additional ports into the data-cache tags when used in a superscalar processor. In such a processor, the cache-tag ports are one of the most critical resources, and increasing their number may result in a physically larger cache and/or a slower cache, or may be infeasible. Stream buffers, on the other hand, sit on the memory side of the data cache and while they must be probed on every data reference, they do not need access to the cache tags. Stream buffers trigger a prefetch based on previous cache misses. Also, they avoid polluting the cache by placing prefetched data in special buffers.

Unlike non-blocking loads and stream buffers, speculative execution is a software-based technique. It involves moving code beyond branches (see [16] for more details). Speculative execution of load instructions is not the same as software prefetching using non-blocking loads; the former involves the movement of *existing* load instructions while the latter involves the insertion of *additional* load instructions. Speculative execution can increase the performance of a program in a number of ways. In superscalar designs, the instruction-level parallelism can be increased if there is a sufficient amount of hardware parallelism available. This increase can also provide more flexibility when scheduling load instructions for a machine with extensive support for non-blocking loads. The result of this flexibility is a better tolerance of the cache-miss latency and a reduction in the average per-miss penalty due to allowing more misses to be simultaneously outstanding.

The three techniques we consider can be used alone or together. To illustrate their use, consider the code segment in Figure 1 which, for simplicity, corresponds to a single-issue machine. Assume that both loads miss in the cache and that the miss penalty is 10 cycles. With none of the three techniques in use, the processor has to stall twice and each time for the full length of the miss penalty. If, however, non-blocking loads are used with a lockup-free cache, the processor will have to stall only once, that is, when it tries to execute instruction 5 (due to a true-data dependency between instructions 5 and 3). If instead stream buffers are used, there will still be two stalls but the cache miss caused by the first load will initiate the prefetch for the data required by the next load (see Section 2.2 for details on how a stream buffer works). Hence, the time to service the second load is reduced. Finally, if both non-blocking

```

L2  subi  r1, 1, r1      ; r1 ← r1-1
    bnz   r1, L1         ; branch if r1 is non zero (1)
    ld    r6, 700(sp)    ; r6 ← [sp+700] (2)
    ld    r8, 732(sp)    ; r8 ← [sp+732] (3)
    addi  sp, 8, sp      ; sp ← sp+8 (4)
    mult  r8, r8, r4     ; r8 ← r8 × r4 (5)
    br    L2            ; branch always

```

Figure 1: Code segment for a RISC single-issue processor.

loads and stream buffers are used together, there will be only one stall and the stall time will be smaller than it was with only non-blocking loads. Stream buffers used alone or together with non-blocking loads also reduce the stall time for subsequent iterations of the loop because the prefetch for the required data will have been initiated in a previous iteration.

Speculative execution is useful for increasing the instruction-level parallelism and for improving the effectiveness of non-blocking loads. The first of these effects cannot be applied to this example because we are assuming a single-issue machine. To show the second one, assume we are using non-blocking loads. With speculative execution, the compiler can move the two load instructions above the branch and thus execute them earlier. The end result is that the processor will stall for less time when it executes instruction 5. The use of stream buffers reduces the stall time further because a prefetch of the data for instruction 3 will be issued before this load is executed.

Previous Work

Other researchers have investigated non-blocking loads, prefetching, or speculative loads, but not their combination. Rogers and Li [13] investigated software support for speculative loads with non-blocking caches using many of the Livermore loop kernels; they compared the results to blocking caches. Sohi and Franklin, on the other hand, studied non-blocking loads [14], while Callahan and Mowry have studied software prefetching for scientific codes [3, 11]. Chen and Baer [4] investigated a combination of non-blocking loads and prefetching, but used a lookahead-PC reference prediction method with a production compiler, instrumented with Pixie, and rescheduled only at a basic block level. Comparing the effectiveness of the different techniques used in these studies is difficult because of different assumptions used by each group of researchers, and the fact that no group looked at all three techniques.

In this paper we look at the relative memory-system performance improvement available from non-blocking loads, hardware prefetching, and speculative execution used individually and in combination. We do this investigation in the context of an advanced quad-issue superscalar machine,

and a pipelined memory system made with recent RAM techniques such as synchronous DRAMs. An important part of this study was the advanced Multiflow Compiler technology [9], which provided trace scheduling and support for speculative loads.

2 Simulation Methodology

We investigated the relative performance of non-blocking loads and stream buffers by examining how their use would affect the performance of current state-of-the-art microprocessor systems. This investigation was carried out by simulating a number of machine configurations using a processor model that resembles a number of commercial processors including the PowerPC 604 [15], the DEC EV5 [5], the MIPS T5 [2] and the SUN Ultrasparc [1]. All configurations used the same hardware for servicing instruction requests and executing instructions. As a result, the contribution to the execution time of a benchmark from instructions was constant for all configurations. The machine configurations differ in the hardware that is available for servicing the processor's requests for data. This fact allows us to better estimate the impact of stream buffers and non-blocking loads.

The processor model implements a RISC processor that can issue four instructions per cycle and uses a conventional, statically scheduled, pipeline. Static scheduling is assumed since the performance benefits of dynamic scheduling are not sufficiently clear in view of the impact the more complex pipeline will have on the cycle time of the processor. For example, the statically scheduled DEC Alpha 21064A achieves a clock frequency of 275MHz in a 0.5um technology, while the dynamically scheduled PowerPC 604 achieves a clock frequency of 100MHz in a 0.5um technology. (There are many other differences between these machines, but in general dynamically scheduled machines announced to date have had significantly slower cycle times than statically scheduled machines in the same technology generation.)

The processor we model supports non-blocking stores and can be configured to support non-blocking loads. There are separate instruction and data caches with the instruction cache having a fixed miss penalty. The data cache is always lockup-free irrespective of whether the processor uses blocking or non-blocking loads¹.

The model also can include one or more stream buffers and these are used to prefetch data. The stream buffers may use either unit strides or dynamically calculated strides.

¹The only difference between these two modes of operation is that with blocking loads, no subsequent instruction can execute until the load is resolved. This restriction allows us to use the same cache model for both modes.

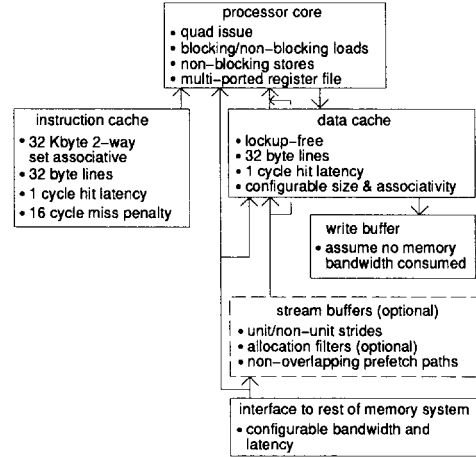


Figure 2: Overview of machine model.

Finally, data is fetched from the next lower level in the memory hierarchy through a bandwidth limited interface. Figure 2 presents an overview of the above machine model. Some of the model details are described further below.

2.1 Processor and Memory Models

Of the four instructions issued per cycle by the processor, each instruction word can contain at most: four integer operations, one floating-point division operation, two floating-point operations, two memory operations (i.e., two loads, two stores, or one of each), and one control flow operation (i.e., branch, subroutine call or return). All integer functional units have single-cycle latencies except for the multiply unit, which is fully pipelined and has a six-cycle latency. All floating point units have three-cycle latencies and are also fully pipelined, with the exception of the floating-point divider. The floating-point divider is not pipelined and has an eight-cycle latency for 32-bit divides, and a 16-cycle latency for 64-bit divides. Finally, stores take one cycle to be resolved and there is a single load-delay slot.

The instruction cache is 32K-byte, 2-way set associative with a 1-cycle cache-hit latency and a 16-cycle cache-miss latency. In addition to cache-miss induced stalls, branch instructions may also introduce stalls if the branch-delay slot(s) cannot be filled by the compiler or if the branch direction is mispredicted. In view of the high correct prediction rates reported by McFarling [10] and Yeh and Patt [17], we assume that all branches are correctly predicted dynamically with the exception of 5% of conditional branches. For conditional branches, the model associates a two cycle

stall with each mispredict. We implement this penalty during the simulation of a benchmark by adding a single-cycle stall for every 10th conditional branch that is executed.

Stores are assumed to be implemented using write-around (i.e., no-write-allocate) and write-through policies with a write buffer situated between the data cache and lower levels in the memory hierarchy. Since our goal is to compare the effectiveness of stream buffers and non-blocking loads while keeping constant other contributions to the execution time, we assume that no memory bandwidth is required to retire stores in the write buffer. This assumption prevents any stalls due to a full write buffer and prevents stores from delaying the servicing of stream buffer or cache fetches.

The lockup-free data cache can resolve cache hits in a single cycle and employs an inverted MSHR (Miss Status Holding Register) organization [6] to process cache misses. An inverted MSHR organization can support as many in-flight cache misses as there are registers and other destinations for data in the processor. Hence, there can be a cache miss outstanding to each of the processor's 32 integer and 32 floating-point registers. The register file has eight read ports and sufficient write ports to prevent any write-port conflicts arising when registers are filled on the resolution of a cache miss.

Requests for blocks of data are sent via the memory interface to the next level in the memory hierarchy. The memory interface returns the requested block in a constant number of cycles, called the *fetch latency*. The bandwidth of the interface is constrained by controlling the number of cycles between the launching of fetch requests. A *fetch spacing* of one allows the memory interface pipeline to be full whereas a spacing equal to the fetch latency allows at most one in-flight fetch. Thus, the time required to resolve a cache miss is not deterministic but has a lower bound equal to the fetch latency. When a block is returned to the cache, the cache line is written simultaneous with the writing of the appropriate words into all registers with loads outstanding to this block (updating all pending registers requires the multiple write ports mentioned above). This simultaneous writing is represented in Figure 2 by the arrows that bypass the data cache. Writing a register or a cache line is assumed to take one cycle.

2.2 Stream Buffers

Our stream buffer model is based on the model originally proposed by Jouppi [8] with the four enhancements described below. In the original model, stream buffers consist of a number of entries that are managed as a FIFO queue. Each entry in the queue can store a block of data and the corresponding address. All entries at the head of the stream-buffer queues are probed at the same time as the

data cache probe is done. If the data cache probe results in a hit, the stream buffers are not touched. However, if a data cache miss occurs, and the desired block is in a head entry, the cache block is read out and it is written into the cache. The stream buffer then issues a prefetch request to the next lower level in the memory hierarchy to fill the empty entry with subsequent blocks.

When a miss occurs to both the cache and the stream buffers, a request for the block containing the miss address is issued. Then, a stream buffer is allocated and told to begin fetching blocks subsequent to the missing block. Because each block that is fetched has a block-address one greater than the last, the stream buffers are said to use a *unit stride*. Once the request for the first block to be prefetched is launched into the memory subsystem, the stream buffer can issue another prefetch request if there remain empty entries in the queue.

Enhancements to the Original Model

In this study we have made four enhancements to the original stream buffer model:

1. allocation filters
2. hardware support for dynamic strides
3. the enforcement of non-overlapping prefetch paths
4. direct access to the stream buffer entries.

These enhancements are described below.

First, in the original stream buffer proposal, a stream buffer is allocated whenever a data reference misses in the cache and in the entries at the head of the stream buffer queues. This allocation policy can result in prefetching down a subsequently unused stream should the data reference be isolated. Prefetching down such a stream will generate excess traffic to the memory system as well as potentially discarding useful data from a previously allocated stream buffer. To prevent this situation from occurring, an *allocation filter* can be used. We implement the filter proposed by Palacharla and Kessler [12]. This filter prevents a stream buffer from being allocated until two misses occur for the same stream. On the second miss, a stream buffer is allocated and it begins prefetching the block subsequent to the one corresponding to the second miss.

Second, instead of limiting prefetch strides to the unit stride of the original proposal, the stride can be determined dynamically based on previous miss addresses. We implemented a scheme based on the *minimum delta* scheme proposed by Palacharla and Kessler [12]. With this scheme, on a stream buffer miss, the allocation filter is applied to determine whether a unit-stride should be used. If there is a filter miss, then the minimum signed difference between the miss address and the last N miss addresses

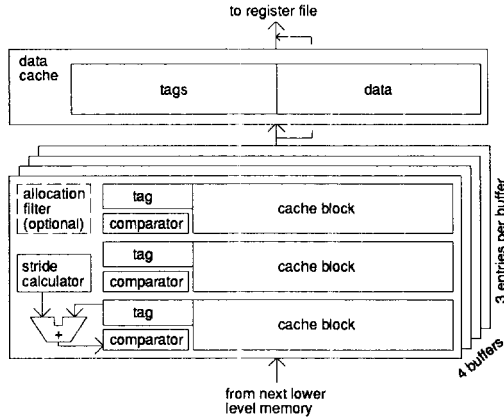


Figure 3: The stream buffer model. In this example, there are four stream buffers with 3 entries each.

is determined; this minimum delta, which may be positive or negative, is the stride. In our model, a stream buffer is allocated if the miss is the third miss in a series to blocks that are separated by this stride. Our stream buffer model with the filter and stride predictor is shown in Figure 3.

Third, when there are multiple stream buffers, we ensure that a given block of data resides in at most one stream buffer. In other words, the stream buffers always prefetch down non-overlapping paths. Non-overlapping paths prevent duplication and thus ensure that the maximum benefit is obtained from the available stream buffers. To achieve non-overlapping paths, a comparator must be associated with each stream buffer entry. While these comparators increase the design complexity, in practical systems, they need to be included for enforcement of multiprocessor cache consistency anyway.

Fourth, we have extended the original stream buffer proposal to include direct access to non-head entries in the queues (this extension is not shown in the figure). This extension reduces the time required to load the cache with data not in the entry at the head of the queue since there is no need to first shift out the blocks closer to the head. We assume that it takes one cycle to extract a block of data from a stream buffer.

2.3 Simulation Framework

To perform the simulations for this study, we used an *object-code translation and instrumentation system*. This system emulates the execution of a benchmark as it would run on a target machine by running the benchmark on an existing machine. As a result *both* the functional behavior and the memory behavior of the application are simulated.

The first step in performing a simulation is to compile the benchmark using instruction scheduling rules pertaining to the architecture of the processor to be modeled. We use a modified version of the Multiflow VLIW Compiler [9] for this purpose². Next, the resulting assembly language (i.e., object code) is translated into the assembly language of the machine on which the simulations are run, namely, Alpha AXP workstations. Instrumentation and modeling code is then inserted into the translated code. Finally, the augmented, translated binary is linked with similarly compiled and instrumented run-time libraries and support routines.

The instrumentation code is inserted to record the emulated run-time behavior of the benchmark. This code records various statistics including cache miss rates, the number of (simulated) instructions executed, and the number of (simulated) clock cycles. The modeling code is inserted to allow the factoring in of the time required to resolve memory and register accesses. This modeling is accomplished by inserting before every emulated load and store instruction a call to a procedure that models the memory. These calls pass to the procedure the address of the item being loaded or stored and the procedure returns the amount of time required to process the access. For example, for non-blocking loads, this time will be the time required to launch the load whereas for a blocking-load it will be the time required to load the data into the cache if it is missing. A mechanism in the simulator adjusts these addresses so that they do not reflect the presence of the simulation infrastructure. Calls to a scoreboard procedure are also inserted before every emulated instruction that uses the result of a load. This procedure factors in the time required to validate the source registers of the instruction.

3 Performance Trends

We investigated four processor designs: one with blocking loads and no stream buffers, known as the *unenhanced design* ("un"), one with non-blocking loads ("nbl"), one with stream buffers ("sb"), and one with both stream buffers and non-blocking loads ("nbl+sb"). For the designs including stream buffers, we considered two types of stream buffers: those that used a unit stride and those that also included an allocation filter and dynamic stride calculator ("+fds"). These six processor design cases are listed in Table 1(a). For those designs including a stream buffer, we assume eight stream buffers each with four entries. The allocation filters and dynamic stride calculators use a table

²The compiler was modified to produce RISC-like object code for a processor with 32-bit addresses, 32-bit integers and 64-bit floating-point numbers. The compiler uses a common backend for both C and Fortran code.

Abbrev.	Processor Details	
	blocking loads	stream buffers
un	yes	none
sb		unit stride
sb+fds		filter & dynamic stride
nbl	no	none
nbl+sb		unit stride
nbl+sb+fds		filter & dynamic stride

(a) Processor Designs

Abbrev.	cache	Memory System Details	
		spacing	latency
ideal	assume 100% data cache hit rate		
8DM	8 KB direct mapped	1	8
64SA	64 KB 2-way associative	8	32

(b) Memory Configurations

Table 1: System details with associated abbreviations.

that stores the 16 last addresses that were found neither in the cache nor in the stream buffers. Since we used 32B cache lines throughout our study, the total data storage of the stream buffers was 1KB.

We studied the relative performance of each of the six processor designs under five memory systems configurations. A memory system configuration comprises a cache and an interface to the next lower level in the memory hierarchy. The organizations we considered are given in Table 1(b). The fetch spacing and latency numbers chosen correspond to a pipelined memory system. These systems are becoming more common with the use of new DRAM technologies such as synchronous DRAMs and pipelined SRAMs. The fetch latency of 8 and spacing of 1 is meant to be representative of microprocessors with two-level on chip caches. Having the backing store for the primary cache on-chip allows a relatively low latency and a very high bandwidth. The fetch latency of 32 and spacing of 8 is intended to be more representative of a system with a single-level of on-chip caching and an optional cache off-chip. For processors with clock frequencies of 200MHz, these latencies and spacings correspond to latencies of 40ns and 160ns and bandwidths of 6.4GB/sec and 800MB/sec. 6.4GB/sec should be achievable on-chip, while 800MB/sec could be easily obtained to an off-chip interface using synchronous DRAMs.

The ideal configuration assumes all data cache references hit in the cache and hence the six processor designs will all achieve the same performance. The other configu-

rations are non-ideal and thus cache misses occur. It is the number of such misses and the time that the processor is stalled that differentiates the six processor designs.

We have simulated the eighteen SPEC92 benchmarks which are listed in Table 2 along with some run-time characteristics. In the table, the columns under the heading "Instructions (millions)" give the dynamic instruction, load and conditional branch counts. Because the same object code and the same input-data sets were used for all simulations of a given benchmark, these numbers remain constant. Thus, we use the average number of cycles per instruction (CPI) as our primary performance measure.

While the numbers of instructions executed are significant, the instruction-word miss rate for each benchmark is usually less than 1%; the exceptions are *doduc* and *xlisp* with a 3% miss rate, and *fpppp* with a 19% miss rate. The next column in the table gives the ideal instructions per cycle, that is, the number of instructions issued per cycle when all stalls are ignored. Observe that the averages are significantly smaller than the maximum of four, a reflection of the scheduling rules and functional unit latencies.

The rest of the columns in the table give statistics for three different memory systems. The first system corresponds to the ideal system and the CPI values for this system are given in the column marked "ideal CPI". In the ideal system, the number of cycles executed for a benchmark is a function of the number of: (1) instruction words executed, (2) stalls caused by functional-unit conflicts, (3) instruction cache misses, and (4) conditional branches. These factors remain constant for all memory configurations and processor designs.

The remaining columns in the figure give statistics for an unenhanced processor using the two memory configurations noted in the caption. For a given benchmark, the portion of the CPI value due to accessing data, the memory CPI, can be found by taking the difference between its ideal CPI and the CPI values measured with a non-ideal memory system. This difference is given in the column with heading "MCPI%" as a percent of the ideal CPI. Observe that for many of the benchmarks the MCPI% value is greater than 20%. Hence, the performance of these application is significantly affected by having to access data, and therefore, it is desirable to reduce the data-access cost.

Note that this way of calculating the MCPI is valid *only* if the processor uses blocking loads. If non-blocking loads are used, a true-data dependency induced stall might be avoided because an instruction cache miss will delay the issuing of the first instruction to use the data. In other words, an instruction cache stall may allow any of the outstanding data cache misses to be resolved, resulting in fewer stall cycles being directly attributable to data references. Thus, for systems with non-blocking loads, it is not possible to

Bench- mark	Instructions (in millions)			Ideal IPC	Ideal CPI	Non-ideal data memory system					
	total	loads	cbranch			8K DM			64K 2-way		
						CPI	MCPI%	ld miss%	CPI	MCPI%	ld miss%
alvinn	3208	942	460	1.85	1.12	1.37	22	10.5	1.72	53	6.4
compress	173	29	16	1.76	0.66	0.96	45	22.3	1.15	74	9.0
dnasa	6858	1644	422	2.16	0.76	1.70	124	49.1	2.21	190	18.9
doduc	1042	238	74	1.90	1.17	1.35	15	10.3	1.20	3	0.5
ear	9506	203	5	1.76	1.03	1.08	5	2.8	1.03	0	≪1
eqntoti	1774	220	189	1.92	0.59	0.64	5	5.5	0.67	14	3.5
espresso	2707	550	402	1.54	0.74	0.91	23	6.5	0.83	12	0.4
fpmp	4294	1131	83	2.45	0.54	1.98	267	11.9	1.74	222	0.1
hydro	5834	1355	328	1.84	0.92	1.28	36	19.1	1.90	107	13.1
mdljdp2	3228	381	313	1.61	1.18	1.34	14	16.7	1.27	8	2.6
mdljsp2	4953	656	88	2.23	0.69	0.76	10	7.2	0.73	5	0.9
ora	4551	90	33	1.68	1.08	1.10	2	1.8	1.08	0	~0
spice	23504	5297	1909	1.33	1.01	1.54	53	29.6	1.65	63	9.0
su2cor	5144	1100	147	2.48	0.50	1.13	56	36.3	1.19	138	10.0
swm	11172	2144	240	2.59	0.47	0.62	32	9.7	1.04	121	9.2
tomcatv	1084	307	14	2.50	0.54	1.10	56	24.9	1.36	152	9.0
wave	3673	694	221	2.14	0.67	0.82	15	8.9	0.77	15	1.4
xlisp	5869	1444	745	1.47	0.83	1.27	53	4.0	1.12	35	0.1

Table 2: Dynamic statistics for each benchmark simulated using an unenhanced processor and three memory system configurations: (1) *ideal*, all data references hit in the cache, (2) an 8K direct mapped cache with a fetch spacing of 1 cycle and a latency of 8, and (3) a 64 KB 2-way set associative cache with a fetch spacing of 8 cycles and a latency of 32. MCPI%, the memory CPI, is the difference between the ideal and non-ideal CPI values as a percent of the ideal CPI.

determine exactly what portion of the CPI is due to accessing data. It is always true, however, that it is better to have a smaller ratio of the non-ideal to ideal CPI values.

We begin by presenting the performance data for *wave* to introduce the our methodology and to point out key characteristics. We then present the data for all benchmarks and discuss the common trends.

3.1 Common Trends

Figure 4(a) presents the CPI for *wave* measured using a memory system with an 8-Kbyte direct-mapped cache, a fetch spacing of 8 and a fetch latency of 32. In this

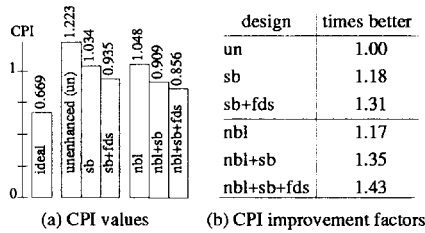


Figure 4: CPI values and ratios for *wave* using an 8 Kbyte direct mapped cache, a fetch spacing of 8 and latency of 32.

figure, each bar corresponds to one of the six processor designs and its height reflects the measured CPI; the numbers above each bar give the actual CPI value. The figure also includes a bar representing the CPI obtained using the ideal memory system. The table in the figure gives the improvement factor for each design in relation to the unenhanced design. These factors are calculated by dividing the unenhanced CPI value by the enhanced CPI value. Note that the CPI factors for the non-blocking load and stream buffer designs cannot be obtained by multiplying together the CPI factor for non-blocking loads and for stream buffers; the combined use of these techniques changes the run-time dynamics of a program.

From the table, we see that the use of unit-stride stream buffers (the “sb” design) results in a CPI improvement of 18%. When allocation filters and dynamic strides are used, this improvement increases to 31%. Non-blocking loads, on the other hand, improve the CPI by 17%, but when used with dynamic strides and allocation filters, the CPI is reduced by 43%. From these percentages, we observe that (1) both non-blocking loads and stream buffers reduce the CPI by a minimum of 17%, (2) non-unit-stride stream buffers give better performance than either non-blocking loads or unit-stride stream buffers, and (3) non-blocking loads and non-unit-stride stream buffers together yield significantly better performance than either technique used alone.

When the miss penalty is reduced, the performance of

		miss penalty	
		large	small
CPI	un	1.233	0.819
times better	sb	1.18	1.06
	sb+fds	1.31	1.09
	nbl	1.17	1.09
	nbl+sb	1.35	1.14
	nbl+sb+fds	1.43	1.15

(a) 8K DM cache CPI improvement factors

		miss penalty	
		large	small
CPI	un	0.767	0.705
times better	sb	1.06	1.01
	sb+fds	1.08	1.02
	nbl	1.01	1.01
	nbl+sb	1.07	1.02
	nbl+sb+fds	1.09	1.02

(b) 64K 2-way SA cache CPI improvement factors

Table 3: CPI values and ratios for *wave* for the large miss penalty (fetch spacing 8 cycles, latency 32 cycles) and the small miss penalty (fetch spacing 1 cycle, latency 8 cycles) memory configurations.

non-blocking loads gets better with respect to the stream-buffer-only designs. This improvement is illustrated by the data presented in Table 3(a). This table gives the unenhanced CPI and the improvement factors for the memory configuration used in Figure 4 and for a memory configuration with a smaller miss penalty. The smaller miss penalty is achieved by decreasing the fetch latency to 8 cycles and increasing the memory bandwidth (by decreasing the fetch spacing to 1 cycle). Observe that the improvement factor for non-blocking loads (“nbl” design) is equal to that for the non-unit-stride stream buffer design (“sb+fds” design) when the miss penalty is smaller. This relative improvement for non-blocking loads occurs with smaller miss penalties because it is more probable that the time required to service a cache miss will be overlapped with the execution of unrelated instructions. Observe also that the improvement factors are smaller. This fact is due to each cache miss requiring less time to be resolved and hence the miss contributes less to the number of cycles executed.

The unenhanced CPI values and improvement factors are given in Table 3(b) for a 64 Kbyte, 2-way set-associative cache and for the same two memory interface configurations. With the larger cache, we see that all the improvement factors have dropped as have the unenhanced CPI

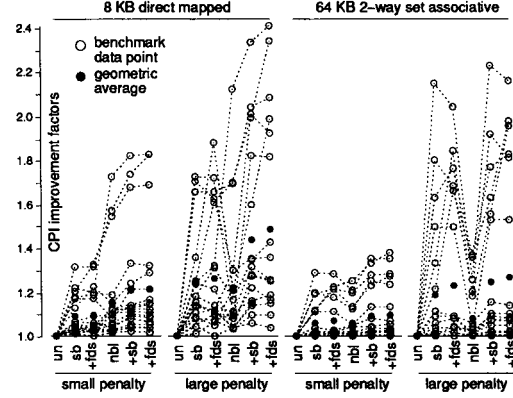


Figure 5: The CPI improvement factors for all benchmarks shown graphically; dotted lines connect points associated with the same benchmark. The “small penalty” points correspond to a fetch spacing of 1 cycle and a fetch latency of 8 while the “large penalty” points correspond to a fetch spacing of 8 and a fetch latency of 32.

values. However, observe that same performance relationships between designs mentioned above also apply here.

The performance relationships just discussed in the context of *wave* occur for many of the other benchmarks. The improvement factors for all 18 benchmarks for the same four memory configurations are presented graphically in Figure 5³. In this graph, there is an unfilled circle for each benchmark for each of the 24 machine configurations, and the filled circles give the geometric average of the improvement factors. Observe the following:

- For each memory configuration, there are a number of benchmarks that incur essentially the same CPI for all processor designs. For these benchmarks, data cache accesses contribute little towards the run-time of the benchmark. However, for other benchmarks, the choice of memory system can make a significant impact.
- Though the improvement factors for all benchmarks are quite different for a particular memory configuration, they vary in similar ways with the different processor designs. These variations are reflected in the geometric averages (the filled circles).
- Stream buffers with dynamic strides and allocation filters have somewhat larger improvement factors (i.e., better performance) than unit-stride stream buffers, although for many benchmarks the additional hardware may not be cost-effective.

³The data from which this graph was prepared is in [7].

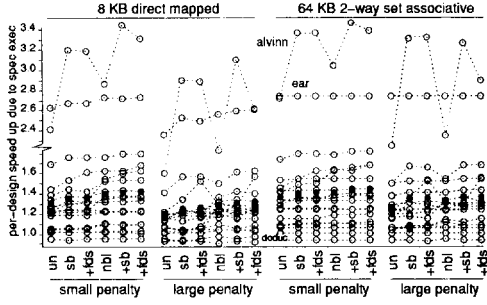


Figure 6: Run-time speedup for each processor design and memory configuration brought about by compiling the code with speculative execution. The “small penalty” points correspond to a fetch spacing of 1 cycle and a fetch latency of 8 while the “large penalty” points correspond to a fetch spacing of 8 and a fetch latency of 32.

- The combined use of stream buffers and non-blocking loads gives a significantly larger performance increase than either technique alone.
- Non-blocking loads have smaller improvement factors (i.e., are less effective) than stream buffers for the large miss penalty. With smaller miss penalties, non-blocking loads are relatively more effective.
- Non-blocking loads are relatively more effective in smaller caches than in larger ones.

Speculative Execution

The effectiveness of non-blocking loads can be increased if more unrelated instructions can be scheduled between the load and first-use instructions. Loop unrolling, and the more general notion of trace scheduling, are techniques that aid in this task by increasing the size of the basic block and hence the pool of unrelated instructions. The results we have presented so far correspond to the use of both of these techniques. Another technique is to allow the compiler to move safe instructions past branch points and thereby allow the instruction to be executed earlier. The Multiflow compiler can implement such code movements by using speculative execution.

To investigate the effect of speculative execution, we re-compiled the benchmarks with speculative execution enabled and simulated their execution on the 24 machine configurations (six processor designs times four memory configurations); speculative execution was allowed for all instructions with the exception of stores and floating-point divides. The two main observations from this investigation are that speculative execution (1) improves by 20 to 40% the

run time of about half the benchmarks on all processor designs, and (2) increases the effectiveness of non-blocking loads. Data supporting the first of these observations is shown graphically in Figure 6; Section 4.2 presents data for the second. This graph presents the speedup in the run-time due to speculative execution. The speedup is calculated on a per-processor-design basis by taking the run-time without speculative execution and dividing it by the run-time with speculative execution. The run-time is used here rather than CPI because speculative execution results in an increase in the number of instructions executed; this increase is discussed further in Section 4.2.

4 Understanding Performance

In the preceding section, we presented data showing that stream buffers and non-blocking loads are effective in improving the performance of a quad-issue processor. We have also seen that speculative execution increases the performance of all six processor designs. In this section, we explore further the causes for the performance gains and explain why these causes are not tied to the specific architecture we simulated.

4.1 Without Speculative Execution

We have seen that non-blocking loads, when acting alone, tend to be more effective when used with caches that have higher miss rates (e.g., smaller caches with less associativity). This is because when the miss density is higher there is a higher probability of being able to overlap more than one miss at a time.

We have also seen that stream buffers, when acting alone, tend to be more effective than non-blocking loads. For non-blocking loads, the fetching of missing data is initiated by the instruction that loads the data into a register. For stream buffers, however, the fetch can be initiated by a previous cache miss with the subsequently executed load instruction loading the register. These fetch-initiating events in effect implement the prefetching of the source operands for a dependent instruction. Non-blocking loads can be viewed as implementing *after-the-load*, or *explicit* prefetching whereas stream buffers can be viewed as implementing *before-the-load*, or *predictive* prefetching. The goal for after-the-load prefetching is to have the data needed by an instruction loaded into registers before the instruction is executed. The goal for before-the-load prefetching is to have the data nearby so that it can be quickly loaded into the registers when the load is executed.

These two types of prefetching are complementary when used together as our results have shown. When used together, the goal is the same as for after-the-load prefetching

but the fetch trigger can be the load or can be a previously detected miss.

Stream buffers work well as a before-the-load prefetcher if there is a reasonable correlation between the block addresses of cache misses since stream buffers must guess what to fetch. If a correct guess is made, the reduction in the stall time is a function of the temporal separation between the fetch trigger and the execution of the load instruction for the data. In our simulations of the SPEC benchmarks with single-, dual- and quad-issue machines, we have observed two trends among the benchmarks that explain why stream buffers work. First, the cache miss addresses are highly predictable and thus it is likely that the correct data will be prefetched. And second, with larger caches, cache misses tend to occur at uneven intervals and in groups. This grouping results in there being time spans when no cache misses occur thereby allowing the stream buffers to fill all their entries. As a result, when a cache miss does occur, it is more likely that the needed data is already resident in a stream buffer. Hence, the stall time is much shorter.

In our simulations of the quad-issue processor, we have observed that most of the SPEC92 benchmarks achieve over a 50% hit rate in the stream buffers, an indication of the sequential nature of the miss addresses. Some of the data on which this observation is based is available in a technical report [7]. We have also found that the unit-stride (32 bytes) is by far the most popular stride. The preference for the unit stride suggests that the allocation filter is more important than the dynamic stride calculator. This result differs from the results reported by Palacharla and Kessler [12] in their study of the NAS and PERFECT benchmarks. They found that some of the benchmarks obtain noteworthy improvement with the use of dynamic strides. That we do not see the same trend in our results is perhaps an artifact of the compiler we are using, or more likely, an artifact of the SPEC92 benchmarks.

The allocation filter prevents a stream buffer from being allocated when a miss occurs to a rarely used part of the address space. By preventing such allocations, prefetched data is likely to remain longer in a stream buffer and thus it is more likely to be used to satisfy a cache miss. For the machine configuration discussed above, the average percent of stream buffer prefetches that are used to resolve a cache miss increases from 22% to 47% when allocation filters are employed. Allocation filters also help by reducing the bandwidth consumed by stream buffer prefetches and thus reduce the time required to fetch data missing from the cache. For some benchmarks (e.g., *alvin* and *tomcatv*), there is very little change in the bandwidth utilization because these benchmarks rarely access infrequently-used parts of their address space. For other benchmarks (e.g.,

compress and *dnasa*), there are many such references and the allocation filters can dramatically decrease the bandwidth utilization. The most dramatic drop occurs for *compress* where the average number of in-progress memory fetches drops from 2 to 0.5.

4.2 With Speculative Execution

With speculative execution we observed that its use (1) improves by 20 to 40% the run time of about half the benchmarks on all processor designs, and (2) increases the effectiveness of non-blocking loads. There are a number of effects that give rise to this behavior. First, the use of speculative execution gives rise to an increase in the average number of instructions issued per cycle (ideal IPC). This fact is shown in Table 4 for each benchmark by the data in the columns headed "ideal IPC". Note that these average values apply to all 24 machine configurations because the number of instructions issued per cycle is determined at compile time. A second effect is that a given benchmark spends less time stalled due to functional-unit conflicts. This reduction occurs because with speculative execution, the compiler can schedule potentially useful instructions between the two that caused the stall rather than scheduling the stall. A third effect is that the compiler can schedule load instructions earlier. This effect is shown in the table by the increase in the weighted average number of instruction issue cycles between the load instruction and the first instruction to use the loaded value (see the columns headed "UAL distance"). Observe that while some benchmarks show greater than a 100% increase, many do not. This lack of change is in part due to our scheduling the code for cache-hit penalty rather than a larger value. The benefits of scheduling for larger values have been demonstrated by Farkas and Jouppi [6] though without the use of speculative execution. However, in quad issue machines scheduling for significantly larger latencies than the cache hit latency becomes infeasible due to a lack of sufficient registers. Due to the small change in the load-use separation, the only configurations that show a significant performance improvement when enhancing speculative execution with non-blocking loads are those using the smaller cache. This result is due to smaller caches being more sensitive to the small changes in non-blocking load dynamics.

5 Conclusions

We have investigated the relative performance impact of non-blocking loads, stream buffers, and speculative loads used individually and in conjunction with each other. We used a quad-issue microprocessor that resembles a number of state-of-the art commercial microprocessors. We have

benchmark	UAL distance		ideal IPC	
	without	with	without	with
alvinn	1.4	3.1	1.8	3.6
compress	1.9	1.9	1.8	2.1
dnasa	4.0	4.5	2.2	3.0
doduc	6.8	6.5	1.9	2.4
ear	2.0	9.7	1.8	3.2
eqntott	1.5	2.2	1.9	3.4
espresso	1.8	2.3	1.5	2.4
fpppp	5.1	4.9	2.4	2.6
hydro	3.3	4.2	1.8	2.5
mdljdp2	4.7	7.6	1.6	2.5
mdljsp2	5.9	6.9	2.2	2.9
ora	5.8	5.2	1.7	1.8
spice	3.0	3.0	1.3	2.1
su2cor	4.8	4.7	2.5	2.9
swm	3.2	5.4	2.6	3.7
tomcatv	5.0	5.1	2.5	2.7
wave	3.2	4.0	2.1	2.7
xlisp	1.8	2.1	1.5	1.9
average	3.2	4.2	1.9	2.6

Table 4: The effects of with and without speculative execution. The first set of columns (UAL distance) gives the weighted average number of instruction issue cycles between a load instruction and the first instruction to use the target register. The second set of columns (ideal IPC) gives the average number of instructions issued per cycle in a machine without stall cycles.

simulated 18 of the SPEC92 benchmarks and evaluated the three techniques by their ability to reduce the number of clock cycles required to execute each benchmark. An important part of this study was the use of the Multiflow Compiler Technology to compile the benchmarks. Using the compiler, we were able to generate object code that was optimized for the microprocessor architecture that we modeled, and to employ trace scheduling and speculative execution, both of which are important for wider-issue machines.

The combined use of stream buffers and non-blocking loads yields significantly better performance than is achieved with either technique acting alone. This complementary behavior is a result of stream buffers being good at reducing the cost of servicing a miss when one occurs, while non-blocking loads are good at hiding the cost of servicing a miss.

Speculative execution was found to improve the performance by 20% to 40% of processors using neither non-blocking loads nor stream buffers as well as those using one or both of these techniques. The performance gains

were found to be fairly constant across all processor designs for a given cache configuration and miss penalty. This performance gain occurred because the compiler was able to schedule approximately 37% more instructions per instruction word and to reduce the number of stalls caused by functional unit conflicts. With speculative execution, the non-blocking load effectiveness increased but this increase was noticeable only with the smaller cache configuration. The non-blocking load improvement occurred because speculative execution resulted in a small increase (from 3.2 instruction issue cycles to 4.2 instruction issue cycles) in the distance between a load instruction and the first instruction to use the loaded value. This increase is nevertheless significant in view of the issue width of the processor we modeled.

The primary benefit from the use of allocation filters and dynamic strides is a reduction in the memory bandwidth consumed by prefetching. In our simulations of the SPEC92 benchmarks, we found that the SPEC92 benchmarks are dominated by unit-stride memory accesses. This observation is in contrast to the conclusion reached by Palacharla and Kessler that some of the NAS and PERFECT benchmarks favor dynamic strides.

Finally, consistent with previous studies, we have found that stream buffers are better able to tolerate larger cache miss penalties, and that the effectiveness of non-blocking loads is improved with smaller caches. This improvement appears to be caused by the higher miss rate and the higher frequency of misses.

Based on our results, the combination of speculative non-blocking loads and stream buffers can reduce the CPI incurred in a blocking system without stream buffers by an average 37% for a 64-Kbyte cache when used with a memory system that can service one request every 8 cycles and has a 32 cycle latency; with an 8-Kbyte cache, the improvement jumps to 61%. We expect that the combination of these techniques will be crucial in removing memory system bottlenecks for processors that attempt to aggressively exploit instruction-level parallelism.

Acknowledgments

The research described in this paper has been partially funded by the National Sciences and Engineering Research Council of Canada and by Digital Equipment Corporation.

We thank Joel Emer, Geoff Lowney, Bob Nix, and David Webb for their guidance and answers to numerous questions as we modified and used the simulation infrastructure they developed. We also thank Alan Eustace, Joel McCormack, Amitabh Srivastava, Annie Warren, and the other WRL-ites for both helping out and putting up with the simulations.

Finally, we thank Digital Equipment Corporation for providing us with the Alpha AXP workstations on which we ran ≈ 6000 simulations, requiring ≈ 540 days of run-time, in 5 months.

References

- [1] Anant Agrawal. Utrasparc: A 64-bit, high-performance sparc processor. *In the proceedings of MicroProcessor Forum*, October 1994.
- [2] John Brennan. T5: A high-performance superscalar mips processor. *In the proceedings of MicroProcessor Forum*, October 1994.
- [3] David Callahan, Ken Kennedy, and Allan Porterfield. Software prefetching. *Proceedings of the 4th ASPLOS Conference*, pages 40–52, 1991.
- [4] Tien-Fu Chen and Jean-Loup Baer. Reducing memory latency via non-blocking and prefetching caches. *Proceedings of the 5th ASPLOS Conference*, pages 51–61, 1992.
- [5] John Edmondson and Paul Rubinfeld. An overview of the 21164 alpha axp microprocessor. *In the proceedings of Hot Chips VI*, August 1994.
- [6] Keith I. Farkas and Norman P. Jouppi. Complexity/performance tradeoffs with non-blocking loads. *Proceedings of the 21st Intl. Symp. on Computer Architecture*, pages 211–222, 1994.
- [7] Keith I. Farkas, Norman P. Jouppi, and Paul Chow. How useful are non-blocking loads, stream buffers and speculative execution in multiple issue processors? *Digital Equipment Corporation Western Research Laboratory Technical Report*, 94/8.
- [8] Norman P. Jouppi. Improving direct-mapped cache performance by the addition of a small fully-associative cache and prefetch buffers. *Proceedings of the 17th Intl. Symp. on Computer Architecture*, pages 364–373, 1990.
- [9] P. Geoffrey Lowney et al. The multiflow trace scheduling compiler. *Journal of Supercomputing*, 7(1-2):51–142, 1993.
- [10] Scott McFarling. Combining branch predictors. *DEC WRL Technical Note TN-36*, 1993.
- [11] Todd C. Mowry, Monica S. Lam, and Anoop Gupta. Design and evaluation of a compiler algorithm for prefetching. *Proceedings of the 5th ASPLOS Conference*, pages 62–73, 1992.
- [12] Subbarao Palacharla and R. E. Kessler. Evaluating stream buffers as a secondary cache replacement. *Proceedings of the 21st Intl. Symp. on Computer Architecture*, pages 24–33, 1994.
- [13] Anne Rogers and Kai Li. Software support for speculative loads. *Proceedings of the 5th ASPLOS Conference*, pages 38–50, 1992.
- [14] Gurindar Sohi and Manoj Franklin. High-bandwidth data memory systems for superscalar processors. *Proceedings of the 4th ASPLOS Conference*, pages 53–62, 1991.
- [15] S. Peter Song. Power pc 604. *In the proceedings of Hot Chips VI*, August 1994.
- [16] M. Srinivas, Alexandru Nicolau, and Vickie H. Allan. An approach to combine predicated/speculative execution for programs with unpredictable branches. *In the proceedings of the International Conference on Parallel Algorithms and Compilation Techniques*, 1994.
- [17] Tse-Yu Yeh and Yale N. Patt. Alternative implementations of two-level adaptive branch prediction. *Proceedings of the 20th Intl. Symp. on Computer Architecture*, pages 124–134, May 1992.