

**Feedback Directed Prefetching: Improving the Performance and
Bandwidth-Efficiency of Hardware Prefetchers**

Santhosh Srinath Onur Mutlu Hyesoon Kim Yale N. Patt



High Performance Systems Group
Department of Electrical and Computer Engineering
The University of Texas at Austin
Austin, Texas 78712-0240

TR-HPS-2006-006
May 2006

This page is intentionally left blank.

Feedback Directed Prefetching: Improving the Performance and Bandwidth-Efficiency of Hardware Prefetchers

Santhosh Srinath Onur Mutlu Hyesoon Kim Yale N. Patt

Department of Electrical and Computer Engineering
The University of Texas at Austin
{santhosh,onur,hyesoon,patt}@ece.utexas.edu

Abstract

High performance processors employ hardware data prefetching to reduce the negative performance impact of large main memory latencies. While prefetching improves performance substantially on many programs, it can significantly reduce performance on others. Also, prefetching can significantly increase memory bandwidth requirements. This paper proposes a mechanism that incorporates dynamic feedback into the design of the prefetcher to increase the average performance improvement provided by prefetching as well as to reduce the negative performance and bandwidth impact of prefetching. Our mechanism estimates prefetcher accuracy, prefetcher timeliness, and prefetcher-caused cache pollution to adjust the aggressiveness of the data prefetcher dynamically. We introduce a new method to track cache pollution caused by the prefetcher at run-time. We also introduce a mechanism that dynamically decides where in the LRU stack to insert the prefetched blocks in the cache based on the cache pollution caused by the prefetcher.

Using the proposed dynamic mechanism improves average performance by 6.5% on the 17 memory-intensive benchmarks in the SPEC CPU2000 suite compared to the best-performing conventional stream-based data prefetcher configuration, while it consumes 18.7% less memory bandwidth. Compared to a conventional stream-based data prefetcher configuration that consumes similar amount of memory bandwidth, feedback directed prefetching provides 13.6% higher performance. Our results show that feedback-directed prefetching eliminates the large negative performance impact incurred on some benchmarks due to prefetching, and it is applicable to stream-based prefetchers, global-history-buffer based delta correlation prefetchers, and PC-based stride prefetchers.

1. Introduction

Hardware data prefetching works by predicting the memory access pattern of the program and speculatively issuing prefetch requests to the predicted memory addresses before the program accesses those addresses. Prefetching has the potential to improve performance if the memory access pattern is correctly predicted and the prefetch requests are initiated early enough before the program accesses the predicted memory addresses. Since the memory latencies faced by today's processors are on the order of hundreds of processor clock cycles, accurate and timely prefetching of data from main memory to the processor caches can lead to significant performance gains by hiding the latency of memory accesses.

On the other hand, prefetching can negatively impact the performance and energy consumption of a processor due to two major reasons, especially if the predicted memory addresses are not accurate:

- First, prefetching can increase the contention for memory bandwidth available in the processor system. The additional bandwidth contention caused by prefetches can lead to increased DRAM bank conflicts, DRAM page conflicts, memory bus contention, and queueing delays. This can significantly reduce performance if it results in delaying demand (i.e.

load/store) requests. Furthermore, inaccurate prefetches can increase the energy consumption of the processor because they result in unnecessary memory accesses (i.e. waste memory/bus bandwidth). The bandwidth contention due to prefetching will become more significant as more and more processing cores are integrated onto the same die in chip multiprocessors, effectively reducing the memory bandwidth available to each core. Therefore, techniques that reduce the memory bandwidth consumption of hardware prefetchers while maintaining their performance improvement will become more desirable and valuable in future processors [20].

- Second, prefetching can cause cache pollution if the prefetched data displaces cache blocks that will later be needed by load/store instructions in the program.¹ Cache pollution due to prefetching may not only reduce performance but also waste memory bandwidth by resulting in additional cache misses.

Furthermore, prefetcher-caused cache pollution generates new cache misses and those generated cache misses in turn generate new prefetch requests. Hence, the prefetcher itself is a *positive feedback system* which can be very unstable in terms of both performance and bandwidth consumption. Therefore, we would like to augment the prefetcher with a *negative feedback system* to make the prefetching system more stable.

Figure 1 compares the performance of varying the aggressiveness of a stream-based hardware data prefetcher from *No prefetching* to *Very Aggressive prefetching* on the 17 memory-intensive benchmarks in the SPEC CPU2000 benchmark suite.² Aggressive prefetching improves IPC performance by 84% on average³ compared to no prefetching and by over 800% for some benchmarks (e.g. mgrid). Furthermore, aggressive prefetching on average performs better than conservative and middle-of-the-road prefetching. Unfortunately, aggressive prefetching significantly reduces performance on some benchmarks. For example, an aggressive prefetcher reduces the IPC performance of ammp by 48% and applu by 29% compared to no prefetching. Hence, blindly increasing the aggressiveness of the hardware prefetcher can drastically reduce performance on several applications even though it improves the average performance of a processor. Since aggressive prefetching significantly degrades performance on some benchmarks, many modern processors employ relatively conservative prefetching mechanisms where the prefetcher does not stay far ahead of the demand access stream of the program [6, 22].

The goal of this paper is to reduce the negative performance and bandwidth impact of aggressive data prefetching while preserving the large performance benefits provided by aggressive prefetching. To achieve this goal, we propose simple and implementable mechanisms that dynamically adjust the aggressiveness of the hardware prefetcher as well as the location in the processor cache where prefetched data is inserted.

The proposed mechanisms estimate the effectiveness of the prefetcher by monitoring the accuracy and timeliness of the prefetch requests as well as the cache pollution caused by the prefetch requests. We describe simple hardware implementations to estimate accuracy, timeliness, and cache pollution. Based on the run-time estimations of these three metrics, the aggressiveness of the hardware prefetcher is decreased or increased dynamically. Also, based on the run-time estimation of the cache pollution caused by the prefetcher, the proposed mechanism dynamically decides where to insert the prefetched blocks in the processor cache's LRU stack.

¹Note that this is a problem only in designs where prefetch requests bring data into processor caches rather than into separate prefetch buffers [12, 10]. In many current processors (e.g. Intel Pentium 4 [6] or IBM POWER4 [22]), prefetch requests bring data into the processor caches. This reduces the complexity of the memory system by eliminating the need to design a separate prefetch buffer. It also makes the large L2 cache space available to prefetch requests, enabling the prefetched blocks and demand-fetched blocks to share the available cache memory dynamically rather than statically partitioning the storage space for demand-fetched and prefetched data.

²Aggressiveness of the prefetcher is determined by how far the prefetcher stays ahead of the demand access stream of the program as well as how many prefetch requests are generated, as shown in Table 1 and described in Section 2.1. Our simulation methodology, benchmarks, and prefetcher configurations are described in detail in Section 4.

³Similar results were reported by [7] and [16]. All average IPC results in this paper are computed as geometric mean of the IPC's of the benchmarks.

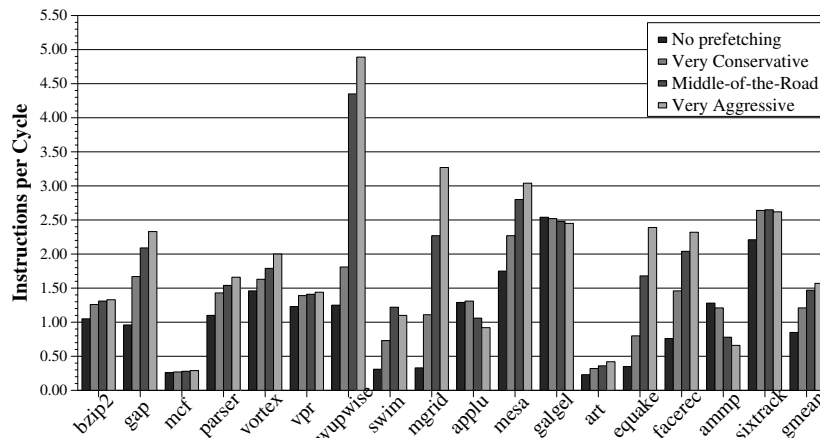


Figure 1. IPC performance as the aggressiveness of a stream-based prefetcher is increased.

Our results show that using the proposed dynamic feedback mechanisms improve the average performance of the 17 memory-intensive benchmarks in the SPEC CPU2000 suite by 6.5% compared to the best-performing conventional stream-based prefetcher configuration. With the proposed mechanism, the negative performance impact incurred on some benchmarks due to stream-based prefetching is completely eliminated. Furthermore, the proposed mechanism consumes 18.7% less memory bandwidth than the best-performing stream-based prefetcher configuration. Compared to a conventional stream-based prefetcher configuration that consumes similar amount of memory bandwidth, feedback directed prefetching provides 13.6% higher performance. We also show that the dynamic feedback mechanism works similarly well when implemented to dynamically adjust the aggressiveness of a global-history-buffer (GHB) based delta correlation prefetcher [9] or a PC-based stride prefetcher [1]. Compared to a conventional GHB-based delta correlation prefetcher configuration that consumes similar amount of memory bandwidth, the feedback directed prefetching mechanism provides 9.9% higher performance. The proposed mechanism provides these benefits with a modest hardware storage cost of 2.54 KB and without significantly increasing hardware complexity. On the remaining 9 SPEC CPU2000 benchmarks, the proposed dynamic feedback mechanism performs as well as the best-performing conventional stream prefetcher configuration for those 9 benchmarks.

We first provide brief background on hardware data prefetching in the next section. Section 3 describes the details and implementation of the proposed feedback-directed prefetching mechanism. Our experimental evaluation methodology is presented in Section 4. We provide the results of our experimental evaluations and our analyses in Section 5. Finally, we describe related research in Section 6.

2. Background and Motivation

2.1. Stream Prefetcher Design

We briefly provide background on how our baseline stream prefetcher works. The stream prefetcher we model is based on the stream prefetcher in the IBM POWER4 processor [22] and more details on the implementation of stream-based prefetching can be found in [10, 17, 22]. The modeled stream prefetcher brings cache blocks from the main memory to the last-level cache, which is the second-level (L2) cache in our baseline processor.

The stream prefetcher is able to keep track of multiple different access streams. For each tracked access stream, a stream tracking entry is created in the stream prefetcher. Each tracking entry can be in one of four different states:

1. Invalid: The tracking entry is not allocated a stream to keep track of. Initially, all tracking entries are in this state.
2. Allocated: A demand (i.e. load/store) L2 miss allocates a tracking entry if the demand miss does not find any existing tracking entry for its cache-block address.
3. Training: The prefetcher trains the direction (ascending or descending direction) of the stream based on the next two L2 misses which occur +/- 16 cache blocks from the first miss.⁴ If the next two accesses in the stream are to ascending (descending) addresses, the direction of the tracking entry is set to 1 (0) and the entry transitions to *Monitor and Request* state.
4. Monitor and Request: The tracking entry monitors the accesses to a memory region from a beginning pointer (address A) to an ending pointer (address P). The maximum distance between the beginning pointer and the ending pointer is determined by *Prefetch Distance*, which indicates how far ahead of the demand access stream the prefetcher can send requests. If there is a demand L2 cache access to a cache block in the monitored memory region, the prefetcher requests cache blocks [P+1, ..., P+N] as prefetch requests (assuming the direction of the tracking entry is set to 1). N is called the *Prefetch Degree*. After sending the prefetch requests, the tracking entry starts monitoring the memory region between addresses A+N to P+N (i.e. effectively it moves the tracked memory region by N cache blocks).⁵

Prefetch Distance and *Prefetch Degree* determine the aggressiveness of the prefetcher. In a traditional prefetcher configuration, the values of *Prefetch Distance* and *Prefetch Degree* are fixed at the design time of the processor. In the feedback directed mechanism we propose, the processor dynamically changes the *Prefetch Distance* and *Prefetch Degree* to adjust the aggressiveness of the prefetcher.

2.2. Metrics of Prefetcher Effectiveness

We use three metrics (*Prefetch Accuracy*, *Prefetch Lateness*, and *Prefetcher-Generated Cache Pollution*) as feedback inputs to feedback directed prefetchers. All three metrics are essential to dynamically decide the aggressiveness of a prefetcher. In this section, we define the metrics and describe the relationship between the metrics and the performance provided by a conventional prefetcher. We evaluate four configurations: *No prefetching*, *Very Conservative prefetching* (distance=4, degree=1), *Middle-of-the-Road prefetching* (distance=16, degree=2), and *Very Aggressive prefetching* (distance=64, degree=4).

2.2.1. Prefetch Accuracy: Prefetch accuracy is a measure of how accurately the prefetcher can predict the memory addresses that will be accessed by the program. Prefetch accuracy is defined using the following equation:

$$\text{Prefetch Accuracy} = \frac{\text{Number of Useful Prefetches}}{\text{Number of Prefetches Sent To Memory}}$$

where *Number of Useful Prefetches* is the number of prefetched cache blocks that are used by demand requests while they are resident in the L2 cache.

Figure 2 shows the IPC (retired Instructions Per Cycle) performance of configurations with different prefetcher aggressiveness along with prefetch accuracy measured over the entire run of each benchmark. The results show that in benchmarks

⁴Note that all addresses tracked by the prefetcher are cache-block addresses.

⁵Right after a tracking entry is trained, the prefetcher sets the beginning pointer to the the first L2 miss address that allocated the tracking entry and the ending pointer to the last L2 miss address that determined the direction of the entry plus an initial start-up distance. Until the monitored memory region's size becomes the same as the *Prefetch Distance* (in terms of cache blocks), the tracking entry increments only the ending pointer by the *Prefetch Degree* when prefetches are issued (i.e. the ending pointer points to the last address requested as a prefetch and the start pointer points to the first L2 miss address that allocated the tracking entry). After the monitored memory region's size becomes the same as *Prefetch Distance*, both the beginning pointer and the ending pointer are incremented by the *Prefetch Degree* (*N*) when prefetches are issued. This way, the prefetcher is able to send prefetch requests that are *Prefetch Distance* ahead of the demand access stream.

where prefetch accuracy is less than 40% (in applu, galgel, and ammp), employing the stream prefetcher always degrades performance compared to the no prefetcher configuration. In all benchmarks where prefetch accuracy is higher than 40% (except mcf), employing the stream prefetcher significantly improves performance over no prefetching. For benchmarks with high prefetch accuracy, performance increases as the aggressiveness of the prefetcher is increased. Hence, the performance improvement provided by increasing the aggressiveness of the prefetcher is correlated with prefetch accuracy.

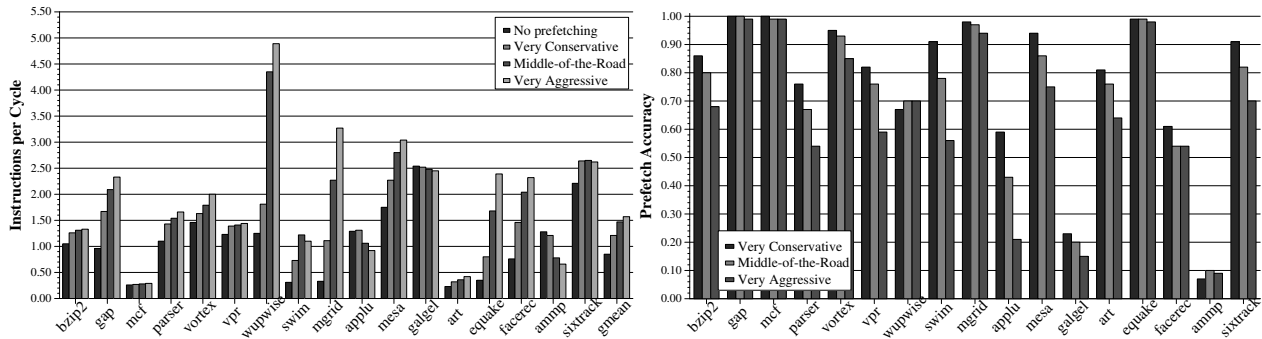


Figure 2. IPC performance (left) and prefetch accuracy (right) with different aggressiveness configurations.

2.2.2. Prefetch Lateness: Prefetch lateness is a measure of how timely the prefetch requests generated by the prefetcher are with respect to the demand accesses that need the prefetched data. A prefetch is defined to be *late* if the prefetched data has not yet returned from main memory by the time a load or store instruction requests the prefetched data. Even though the prefetch requests are accurate, a prefetcher might not be able to improve performance if the prefetch requests are very late. We define prefetch lateness as:

$$Prefetch\ Lateness = \frac{Number\ of\ Late\ Prefetches}{Number\ of\ Useful\ Prefetches}$$

Figure 3 shows the performance of configurations with different prefetcher aggressiveness along with prefetch lateness measured over the entire run of each program. Prefetch lateness results explain why prefetching does not provide significant performance benefit on mcf, even though the prefetch accuracy is close to 100%. More than 90% of the useful prefetch requests are late in mcf. In general, prefetch lateness decreases as the prefetcher becomes more aggressive. For example, in vortex, prefetch lateness decreases from 70% to 22% when a very aggressive prefetcher is used instead of a very conservative one. Aggressive prefetching reduces the lateness of prefetches because an aggressive prefetcher generates prefetch requests earlier than a conservative prefetcher would.

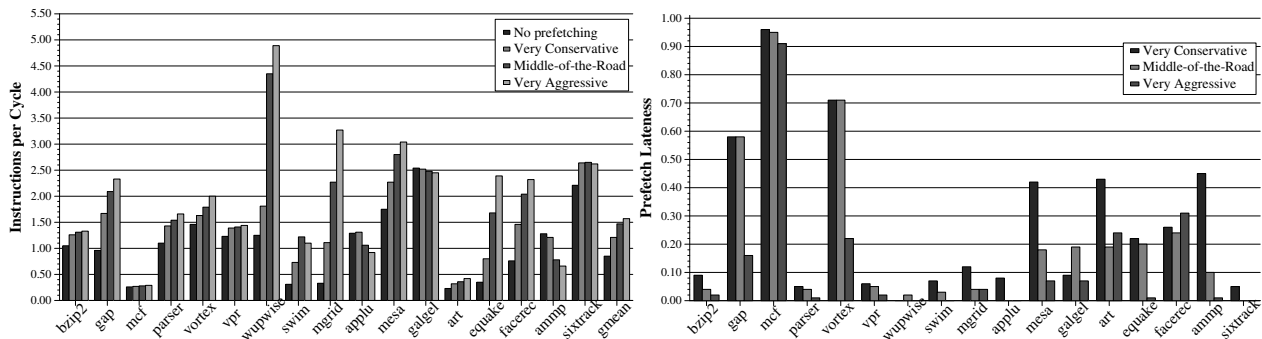


Figure 3. IPC performance (left) and prefetch lateness (right) with different aggressiveness configurations.

2.2.3. Prefetcher-Generated Cache Pollution: Prefetcher-generated cache pollution is a measure of the disturbance caused by prefetched data in the L2 cache. Cache pollution caused by the prefetcher is quantitatively defined as:

$$\text{PrefetcherGenerated Cache Pollution} = \frac{\text{Number of Demand Misses Caused By the Prefetcher}}{\text{Number of Demand Misses}}$$

A demand miss is defined to be caused by the prefetcher if it would not have occurred had the prefetcher not been present. A higher value for *Prefetcher-Generated Cache Pollution* indicates that the prefetcher pollutes the cache more. If the prefetcher-generated cache pollution is high, the performance of the processor can degrade because useful data in the cache could be evicted by prefetched data. Furthermore, high cache pollution can also result in higher memory bandwidth consumption by requiring the re-fetch of the displaced data from main memory.

3. Feedback Directed Prefetching

Feedback directed prefetching dynamically adapts the aggressiveness of the prefetcher based on the accuracy, lateness, and pollution metrics defined in the previous section. In this section, we describe the hardware mechanisms that track these metrics. We also describe our feedback directed prefetching mechanism that uses these metrics to adjust the behavior of the prefetcher.

3.1. Hardware Mechanisms for Collecting Feedback Information

3.1.1. Prefetch Accuracy: To track the usefulness of prefetch requests we add a bit, called the *pref-bit*, to each tag-store entry in the L2 cache.⁶ When a prefetched block is inserted into the cache, the *pref-bit* associated with that block is set. The accuracy of the prefetcher is tracked using two hardware counters. The first counter, *pref-total*, tracks the number of prefetches sent to memory. The second counter, *used-total*, tracks the number of useful prefetches. When a prefetch request is sent to memory, the *pref-total* counter is incremented. When an L2 cache block that has the *pref-bit* set is accessed by a demand request, the *pref-bit* is reset and the *used-total* counter is incremented. The accuracy of the prefetcher is computed by taking the ratio of *used-total* to *pref-total*.

3.1.2. Prefetch Lateness: Miss Status Holding Register (MSHR) [11] is a hardware structure that keeps track of all the in-flight memory requests. Before allocating an MSHR entry for a request, the MSHR checks if the requested cache block is being serviced by an earlier memory request. Each entry in the L2 cache MSHR has a bit, called the *pref-bit*, which indicates that the memory request was generated by the prefetcher. A prefetch request is late if a demand request for the prefetched address is generated while the prefetch request is in the MSHR waiting for main memory. We use a hardware counter, *late-total*, to keep track of such late prefetches. If a demand request hits an MSHR entry that has its *pref-bit* set, the *late-total* counter is incremented, and the *pref-bit* associated with that MSHR entry is reset. The lateness metric is computed by taking the ratio of the number of late prefetch requests (*late-total*) to the total number of useful prefetch requests (*used-total*).

3.1.3. Prefetcher-Generated Cache Pollution: To track the number of demand misses caused by the prefetcher, the processor needs to store information about all the demand-fetched L2 cache blocks dislodged by the prefetcher. However, such a mechanism is impractical as it incurs a heavy overhead in terms of both hardware and complexity. We use the Bloom filter

⁶Note that several proposed prefetching implementations, such as tagged next-sequential prefetching [5, 19] already employ *pref-bits* in the cache.

concept [2, 18] to provide a simple cost-effective hardware mechanism that can approximate the number of demand misses caused by the prefetcher.

Figure 4 shows the filter that is used to approximate the number of L2 demand misses caused by the prefetcher. The filter consists of a bit-vector, which is indexed with the output of the exclusive-or operation of the lower and higher order bits of the cache block address. When a block that was brought into the cache due to a demand miss is evicted from the cache due to a prefetch request, then the filter is accessed with the address of the evicted cache block and the corresponding bit in the filter is set (indicating that the evicted cache block was evicted due to a prefetch request). When a prefetch request is serviced from memory, the pollution filter is accessed with the cache-block address of the prefetch request and the corresponding bit in the filter is reset. When a demand access misses in the cache, the filter is accessed using the cache-block address of the demand request. If the corresponding bit in the filter is set, it is an indication that the demand miss was caused by the prefetcher. In such cases, the hardware counter, *pollution-total*, that keeps track of the total number of demand misses caused by the prefetcher is incremented. Another counter, *demand-total*, keeps track of the total number of demand misses generated by the processor and is incremented for each demand miss. Cache pollution caused by the prefetcher can be computed by taking the ratio of *pollution-total* to *demand-total*. We use a 4096-entry bit vector in all our experiments.

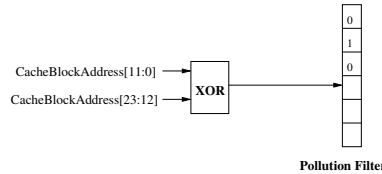


Figure 4. Filter used to determine cache pollution caused by the prefetcher.

3.2. Sampling-based Update of Feedback Information

To adapt to the time-varying memory phase behavior of a program, we use interval-based sampling for all the counters described in Section 3.1. The execution time of a program is divided into intervals and the value of each counter is computed as:

$$CounterValue = \frac{1}{2}CounterValueAtTheBeginningOfTheInterval + \frac{1}{2}CounterValueDuringInterval \quad (1)$$

The *CounterValueDuringInterval* is reset at the end of each sampling interval. The above equation used to update the counters (Equation 1) gives more weight to the behavior of the program in the most recent interval while taking into account the behavior in all previous intervals. Our mechanism defines the length of an interval based on the number of useful cache blocks evicted from the L2 cache.⁷ A hardware counter, *eviction-count*, keeps track of the number of blocks evicted from the L2 cache. When the value of the counter exceeds a statically-set threshold $T_{interval}$, the interval ends. At the end of an interval, all counters described in Section 3.1 are updated according to Equation 1. The updated counter values are then used to compute the three metrics: accuracy, lateness, and pollution. These metrics are used to adjust the prefetcher behavior for the next interval. The eviction-count register is reset and a new interval begins. In our experiments, we use a value of 8192 (half the number of blocks in the L2 cache) for $T_{interval}$.

⁷There are other ways to define the length of an interval, e.g. based on the number of instructions executed. We use the number of useful cache blocks evicted to define an interval because this metric provides a more accurate view of the memory behavior of a program than the number of instructions executed.

3.3. Dynamically Adjusting Prefetcher Behavior

At the end of each sampling interval, the computed values of the accuracy, lateness, and pollution metrics are used to dynamically adjust the behavior of the prefetcher. The behavior of the prefetcher is adjusted in two ways: (1) by adjusting the aggressiveness of the prefetching mechanism, (2) by adjusting the location in the L2 cache’s LRU stack where prefetched blocks are inserted.

3.3.1. Adjusting Prefetcher Aggressiveness: The aggressiveness of the prefetcher directly determines the potential for benefit as well as harm that is caused by the prefetcher. By dynamically adapting this parameter based on the collected feedback information, the processor can not only achieve the performance benefits of aggressive prefetching during program phases where aggressive prefetching performs well but also eliminate the negative performance and bandwidth impact of aggressive prefetching during phases where aggressive prefetching performs poorly.

As shown in Table 1, our baseline stream prefetcher has five different configurations ranging from *Very Conservative* to *Very Aggressive*. The aggressiveness of the stream prefetcher is determined by the *Dynamic Configuration Counter*, which is a 3-bit saturating counter that saturates at values 1 and 5. Initially, the value of the *Dynamic Configuration Counter* is set to 3, which indicates Middle-of-the-Road aggressiveness striking a middle ground between too conservative and too aggressive prefetching.

<i>Dynamic Configuration Counter</i>	<i>Aggressiveness</i>	<i>Prefetch Distance</i>	<i>Prefetch Degree</i>
1	Very Conservative	4	1
2	Conservative	8	1
3	Middle-of-the-Road	16	2
4	Aggressive	32	4
5	Very Aggressive	64	4

Table 1. Stream prefetcher configurations.

At the end of each sampling interval, the value of the *Dynamic Configuration Counter* is updated based on the computed values of the accuracy, lateness, and pollution metrics. The computed accuracy is compared to two thresholds (A_{high} and A_{low}) and is classified as high, medium or low. Similarly, the computed lateness is compared to a single threshold ($T_{lateness}$) and is classified as either late or not-late. Finally, the computed pollution caused by the prefetcher is compared to a single threshold ($T_{pollution}$) and is classified as high (polluting) or low (not-polluting). We use static thresholds in our mechanisms. The effectiveness of our mechanism can be improved by dynamically tuning the values of these thresholds and/or using more thresholds, but such optimization is out of the scope of this paper. In Section 5, we show that even with untuned threshold values, the proposed feedback directed prefetching mechanism can significantly improve performance and reduce memory bandwidth consumption on different data prefetchers.

Table 2 shows in detail how the estimated values of the three metrics are used to adjust the dynamic configuration of the prefetcher. We determined the counter update choice for each case empirically. If the prefetches are causing pollution (all even-numbered cases), the prefetcher is adjusted to be less aggressive to reduce cache pollution and to save memory bandwidth (except in Case 2 when the accuracy is high and prefetches are late – we do increase aggressiveness in this case to gain more benefit from highly-accurate prefetches). If the prefetches are late but not polluting (Cases 1,5,9), then the aggressiveness is increased to increase timeliness unless the prefetch accuracy is low (Case 9 – we do not increase aggressiveness in this case because a large fraction of inaccurate prefetches will waste memory bandwidth). If the prefetches

are neither late nor polluting (Cases 3, 7, 11), the aggressiveness is left unchanged except when the accuracy is low (Case 11 – we reduce aggressiveness in this case to save wasted memory bandwidth due to inaccurate prefetches).

Case	Prefetch Accuracy	Prefetch Lateness	Cache Pollution	Dynamic Configuration Counter Update (reason)
1	High	Late	Not-Polluting	Increment (to increase timeliness)
2	High	Late	Polluting	Increment (to increase timeliness)
3	High	Not-Late	Not-Polluting	No Change (best case configuration)
4	High	Not-Late	Polluting	Decrement (to reduce pollution)
5	Medium	Late	Not-Polluting	Increment (to increase timeliness)
6	Medium	Late	Polluting	Decrement (to reduce pollution)
7	Medium	Not-Late	Not-Polluting	No Change (to keep the benefits of timely prefetches)
8	Medium	Not-Late	Polluting	Decrement (to reduce pollution)
9	Low	Late	Not-Polluting	No Change (to keep the benefits of timely prefetches)
10	Low	Late	Polluting	Decrement (to reduce pollution)
11	Low	Not-Late	Not-Polluting	Decrement (to save bandwidth)
12	Low	Not-Late	Polluting	Decrement (to reduce pollution and save bandwidth)

Table 2. How to adapt? Use of the three metrics to adjust the aggressiveness of the prefetcher.

3.3.2. Adjusting Cache Insertion Policy of Prefetched Blocks: Our feedback directed prefetching mechanism also adjusts the location in which a prefetched block is inserted in the LRU-stack of the corresponding cache set based on the observed behavior of the prefetcher. In many cache implementations, prefetched cache blocks are simply inserted into the Most-Recently-Used (MRU) position in the LRU-stack, since such an insertion policy does not require any changes to the cache implementation. Inserting the prefetched blocks into the MRU position can allow the prefetcher to be more aggressive and request data long before its use because this insertion policy allows the useful prefetched blocks to stay longer in the cache. However, if the prefetched cache blocks create cache pollution, having a different cache insertion policy for prefetched cache blocks can help reduce the cache pollution caused by the prefetcher. A prefetched block that is not useful creates more pollution in the cache if it is inserted into the MRU position rather than a less recently used position because it stays in the cache for a longer time period, occupying cache space that could otherwise be allocated to a useful demand-fetched cache block. Therefore, if the prefetch requests are causing cache pollution, it would be desirable to reduce this pollution by changing the location in the LRU stack in which prefetched blocks are inserted.

We propose a simple heuristic that decides where in the LRU stack of the L2 cache set a prefetched cache block is inserted based on the estimated prefetcher-generated cache pollution. At the end of a sampling interval, the estimated cache pollution metric is compared to two thresholds (P_{low} and P_{high}) to determine whether the pollution caused by the prefetcher was low, medium, or high. If the pollution caused by the prefetcher was low, then the prefetched cache blocks are inserted into the middle (MID) position in the LRU stack during the next sampling interval (for an n-way set-associative cache, we define the MID position in the LRU stack as the $\text{floor}(n/2)$ th least-recently-used position).⁸ On the other hand, if the pollution caused by the prefetcher was medium, the prefetched cache blocks are inserted into the LRU-4 position in the LRU stack (for an n-way set-associative cache, we define the LRU-4 position in the LRU stack as the $\text{floor}(n/4)$ th least-recently-used position). Finally, if the pollution caused by the prefetcher was high, the prefetched cache blocks are inserted into the LRU position during the next sampling interval in order to reduce the pollution caused by the prefetched cache blocks.

⁸We found that inserting prefetched blocks into the MRU position does not provide significant benefits over inserting them in the MID position. Therefore, our dynamic mechanism does not insert prefetched blocks into the MRU position. For a detailed analysis of the cache insertion policy of prefetched blocks, see Section 5.2.

4. Evaluation Methodology

We evaluate the performance impact of feedback directed prefetching on an in-house execution-driven Alpha ISA simulator that models an aggressive superscalar, out-of-order execution processor. The parameters of the processor we model are shown in Table 3. As the performance impacts of prefetching and our mechanisms are dependent on the memory system modeled, we briefly describe our memory model next.

Pipeline	20-cycle minimum branch misprediction penalty; 4 GHz processor
Branch Predictor	aggressive hybrid branch predictor (64K-entry gshare, 64K-entry per-address w/ 64K-entry selector) wrong-path execution faithfully modeled
Instruction Window	128-entry reorder buffer; 128-entry INT, 128-entry FP physical register files; 64-entry store buffer;
Execution Core	8-wide, fully-pipelined except for FP divide; full bypass network
On-chip Caches	64KB Instruction cache with 2-cycle latency; 64KB, 4-way L1 data cache with 8 banks and 2-cycle latency, allows 4 load accesses per cycle; 1MB, 16-way, unified L2 cache with 8 banks and 10-cycle latency, 128 L2 MSHRs, 1 L2 read port, 1 L2 write port; all caches use LRU replacement and have 64B block size
Buses and Memory	500-cycle minimum main memory latency; 32 DRAM banks; 32B-wide, split-transaction core-to-memory bus at 4:1 frequency ratio; 4.5 GB/s bus bandwidth; max. 128 outstanding misses to main memory; bank conflicts, bandwidth, port contention, and queueing delays faithfully modeled

Table 3. Baseline processor configuration.

4.1. Memory Model

We evaluate our mechanisms using a detailed memory model which mimics the behavior and the bandwidth/port limitations of all the hardware structures in the memory system faithfully. All the mentioned effects are modeled correctly and bandwidth limitations are enforced in our model as described in [15]. The major components in the evaluated memory model are shown in Figure 5. The memory bus has a bandwidth of 4.5 GB/s.

The baseline hardware data prefetcher we model is a stream prefetcher that can track 64 different streams. Prefetch requests generated by the stream prefetcher are inserted into the Prefetch Request Queue which has 128 entries in our model. Requests are drained from this queue and inserted into the L2 Request Queue and are given the lowest priority so that they do not delay demand load/store requests. Requests that miss in the L2 cache access DRAM memory by going through the Bus Request Queue. L2 Request Queue, Bus Request Queue, and L2 Fill Queue have 128 entries each. Only when a prefetch request goes out on the bus does it count towards the number of prefetches sent to memory. A prefetched cache block is placed into the MRU position in the L2 cache in the baseline configuration.

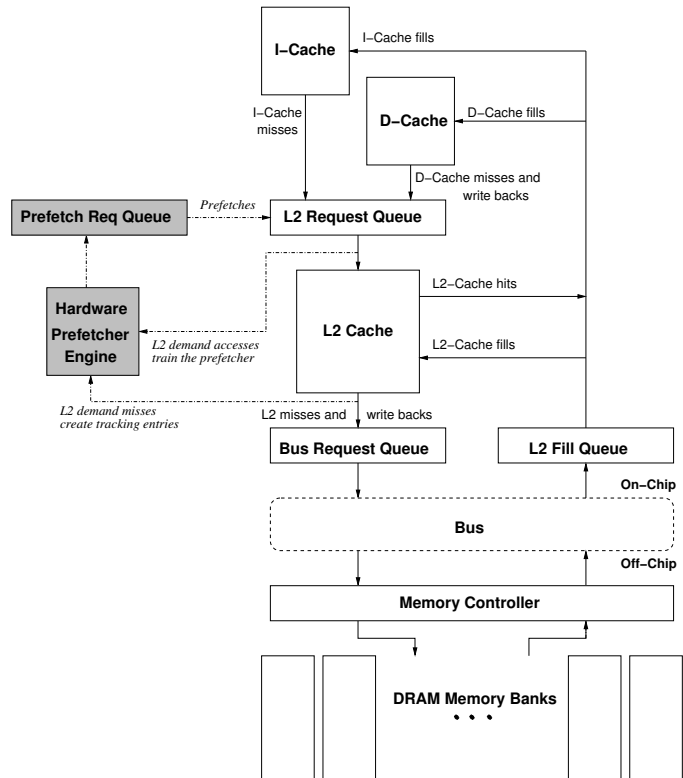


Figure 5. Memory system modeled in our evaluations.

4.2. Benchmarks

We focus our evaluation on those benchmarks from the SPEC CPU2000 suite where the most aggressive prefetcher configuration sends out to memory at least 200K prefetch requests over the 250 million instruction run. On the remaining

nine programs of the SPEC CPU2000 suite, the potential for improving either performance or bandwidth-efficiency of the prefetcher is limited because the prefetcher is not active (even if it is configured very aggressively).⁹ For reference, the number of prefetches generated for each benchmark in the SPEC CPU2000 suite is shown in Table 4.

The benchmarks were compiled using the Compaq C/Fortran compilers with the `-fast` optimizations and profile-driven feedback enabled. All benchmarks are fast forwarded to skip the initialization portion of each benchmark and then simulated for 250 million instructions.

gzip	vpr	gcc	mcf	crafty	parser	eon	perlbnk	gap	vortex	bzip2	twolf
31K	246K	110K	2585K	59K	515K	4969	9218	1656K	591K	336K	2749

wupwise	swim	mgrid	applu	mesa	galgel	art	equake	facerec	ampp	lucas	fma3d	sixtrack	apsi
799K	8766K	2185K	6038K	273K	243K	13319K	2414K	2437K	1157K	1103	3643	292K	8656

Table 4. Number of prefetches sent by a very aggressive stream prefetcher for each benchmark in the SPEC CPU2000 suite. The benchmarks shown in bold are used in our evaluations.

4.3. Thresholds Used in Feedback Directed Prefetching Implementation

The thresholds used in the implementation of our mechanism are provided in Table 5. We determined the parameters of our mechanism empirically. However, we did not tune the parameters to our application set since this requires an exponential number of simulations in terms of the different parameter combinations. We estimate that optimizing these thresholds can further improve the performance and bandwidth-efficiency of our mechanism.

In systems where bandwidth contention is estimated to be higher (e.g. systems where many threads share the memory bandwidth), A_{high} and A_{low} thresholds can be increased to restrict the prefetcher from being too aggressive. In systems where the lateness of prefetches is estimated to be higher due to higher contention in the memory system, reducing the $T_{lateness}$ threshold can increase performance by increasing the timeliness of the prefetcher. Reducing $T_{pollution}$, P_{high} or P_{low} thresholds results in reducing the prefetcher-generated cache pollution. In systems with higher contention for the L2 cache space (e.g. systems with a smaller L2 cache or with many threads sharing the same L2 cache), reducing the values of $T_{pollution}$, P_{high} or P_{low} may be desirable to reduce the cache pollution due to prefetching.

A_{high}	A_{low}	$T_{lateness}$	$T_{pollution}$	P_{high}	P_{low}
0.75	0.40	0.01	0.005	0.25	0.005

Table 5. Thresholds used in the evaluation of feedback directed prefetching.

5. Experimental Results and Analysis

5.1. Adjusting Prefetcher Aggressiveness

We first evaluate the performance of dynamic feedback directed prefetching to adjust the aggressiveness of the stream prefetcher (as described in Section 3.3.1) in comparison to four traditional configurations that do not incorporate dynamic feedback: No prefetching, Very Conservative prefetching, Middle-of-the-Road prefetching, and Very Aggressive prefetching.

⁹We also evaluated the remaining benchmarks that have less potential. Results for these benchmarks are shown in Section 5.11.

Figure 6 shows the IPC performance of each configuration. Adjusting the prefetcher aggressiveness dynamically (i.e. *Dynamic Aggressiveness*) provides the best average performance across all configurations. In particular, dynamically adapting the aggressiveness of the prefetcher using the proposed feedback mechanism provides 4.7% higher average IPC over the Very Aggressive configuration and 11.9% higher IPC over the Middle-of-the-Road configuration.

On almost all benchmarks, *Dynamic Aggressiveness* provides performance that is very close to the performance achieved by the best-performing traditional prefetcher configuration for each benchmark. Hence, the dynamic mechanism is able to detect and employ the best-performing aggressiveness level for the stream prefetcher on a per-benchmark basis.

Figure 6 shows that *Dynamic Aggressiveness* almost completely eliminates the large performance degradation incurred on some benchmarks due to Very Aggressive prefetching. While the most aggressive traditional prefetcher configuration provides the best average performance, it results in a 28.9% performance loss on applu and a 48.2% performance loss on ammp compared to no prefetching. In contrast, *Dynamic Aggressiveness* results in a 1.8% performance improvement on applu and only a 5.9% performance loss on ammp compared to no prefetching, similar to the best-performing traditional prefetcher configuration for the two benchmarks.

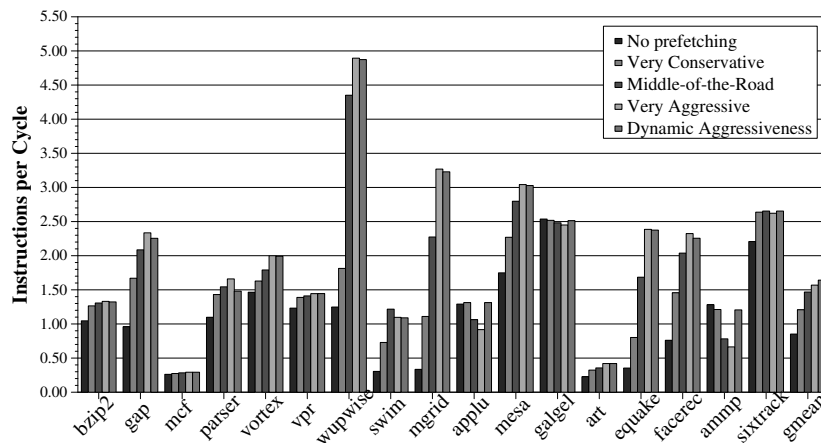


Figure 6. Performance of feedback-directed adjustment of prefetcher aggressiveness in comparison to conventional prefetcher configurations.

5.1.1. Adapting to the Program Figure 7 shows the distribution of the value of the *Dynamic Configuration Counter* over all sampling intervals in the *Dynamic Aggressiveness* mechanism. For benchmarks where aggressive prefetching hurts performance (e.g. applu, galgel, ammp), the feedback mechanism chooses and employs the least aggressive dynamic configuration (counter value of 1) for most of the sampling intervals. For example, the prefetcher is configured to be Very Conservative in more than 98% of the sampling intervals for both applu and ammp.

On the other hand, for benchmarks where aggressive prefetching significantly increases performance (e.g. wupwise, mgrid, equake), feedback directed prefetching employs the most aggressive configuration (counter value of 5) for most of the sampling intervals. For example, the prefetcher is configured to be Very Aggressive in more than 98% of the sampling intervals for wupwise, mgrid, and equake.

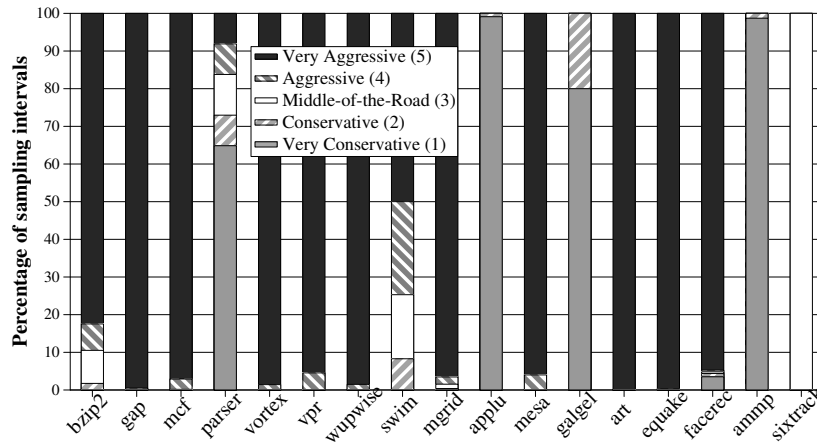


Figure 7. Distribution of the dynamic aggressiveness level (dynamic configuration counter) of the prefetcher in feedback directed prefetching.

5.2. Adjusting Cache Insertion Policy of Prefetched Blocks

Figure 8 shows the performance of dynamically adjusting the cache insertion policy (i.e. *Dynamic Insertion*) using feedback directed prefetching as described in Section 3.3.2. The performance of *Dynamic Insertion* is compared to four different static insertion policies that always insert a prefetched block into the (1) MRU position, (2) MID ($\lfloor n/2 \rfloor$ th) position where n is the set-associativity, (3) LRU-4 ($\lfloor n/4 \rfloor$ th least-recently-used) position, and (4) LRU position in the LRU stack. The dynamic cache insertion policy is evaluated using the Very Aggressive prefetcher configuration.

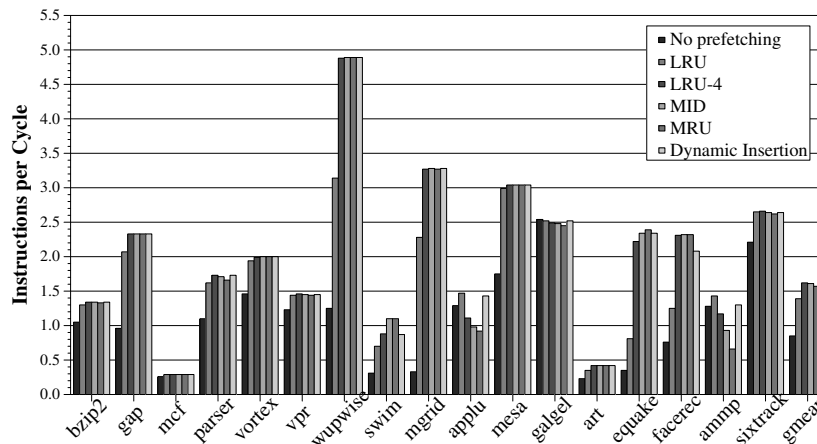


Figure 8. Performance of feedback-directed adjustment of prefetch insertion policy in comparison to static policies.

The data in Figure 8 shows that statically inserting prefetches in the LRU position can result in significant average performance loss compared to statically inserting prefetches in the MRU position. This is because inserting prefetched blocks in the LRU position causes an aggressive prefetcher to evict prefetched blocks before they get used by demand loads/stores. However, inserting in the LRU position eliminates the performance loss due to aggressive prefetching in benchmarks where aggressive prefetching hurts performance (e.g. *applu* and *ammp*). Among the static cache insertion policies, inserting the prefetched blocks into the LRU-4 position provides the best average performance, improving performance by 3.2% over inserting prefetched blocks in the MRU position.

Adjusting the cache insertion policy dynamically provides higher performance than any of the static insertion policies. *Dynamic Insertion* achieves 5.1% better performance than inserting prefetched blocks into the MRU position and 1.9% better performance than inserting them into the LRU-4 position. Furthermore, *Dynamic Insertion* almost always provides the performance of the best static insertion policy for each benchmark. Hence, dynamically adapting the prefetch insertion policy using run-time estimates of prefetcher-generated cache pollution is able to detect and employ the best-performing cache insertion policy for the stream prefetcher on a per-benchmark basis.

Figure 9 shows the distribution of the insertion position of the prefetched blocks when *Dynamic Insertion* is used. For benchmarks where a static policy of inserting prefetched blocks into the LRU position provides the best performance across all static configurations (applu, galgel, ammp), *Dynamic Insertion* places most (more than 50%) of the prefetched blocks into the LRU position. Therefore, *Dynamic Insertion* improves the performance of these benchmarks by dynamically employing the best-performing insertion policy.

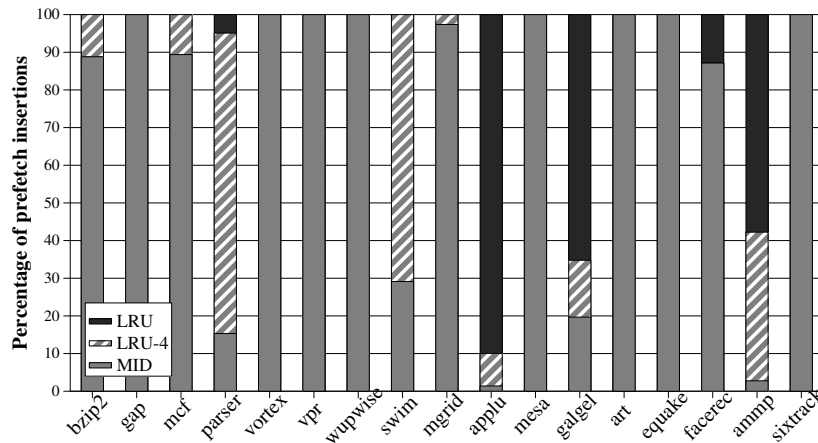


Figure 9. Distribution of the insertion position of prefetched blocks using the feedback-directed cache insertion policy.

5.3. Putting It All Together: Dynamically Adjusting Both Prefetcher Aggressiveness and Insertion Policy

This section examines the use of feedback directed prefetching for dynamically adjusting both the prefetcher aggressiveness (*Dynamic Aggressiveness*) and the cache insertion policy of prefetched blocks (*Dynamic Insertion*). Figure 10 compares the performance of five different mechanisms from left to right: (1) No prefetching, (2) Very Aggressive prefetching, (3) Very Aggressive prefetching with *Dynamic Insertion*, (4) *Dynamic Aggressiveness*, and (5) *Dynamic Aggressiveness* and *Dynamic Insertion* together.

Using *Dynamic Aggressiveness* and *Dynamic Insertion* together provides the best performance across all configurations, improving the IPC by 6.5% over the best-performing traditional prefetcher configuration (i.e. Very Aggressive configuration). This performance improvement is greater than the performance improvement provided by *Dynamic Aggressiveness* or *Dynamic Insertion* alone. Hence, dynamically adjusting both aspects of prefetcher behavior (aggressiveness and insertion policy) provides complementary performance benefits.

With the use of feedback directed prefetching to dynamically adjust both aspects of prefetcher behavior, the performance loss incurred on some benchmarks due to aggressive prefetching is completely eliminated. No benchmark loses performance compared to employing no prefetching if both *Dynamic Aggressiveness* and *Dynamic Insertion* are used. In fact, feedback

directed prefetching improves the performance of applu by 13.4% and ammp by 11.4% over no prefetching – two benchmarks that otherwise incur very significant performance losses with an aggressive traditional prefetcher configuration.

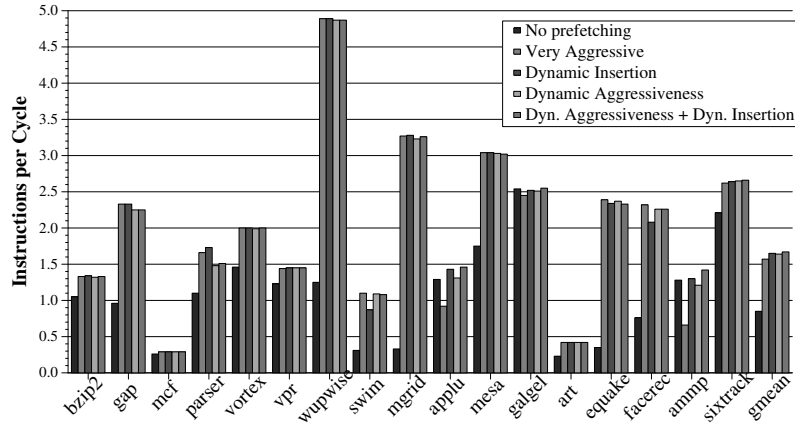


Figure 10. Performance of feedback directed prefetching (adjusting both aggressiveness and insertion policy).

5.4. Impact of Feedback Directed Prefetching on Memory Bandwidth Consumption

Aggressive prefetching can adversely affect the bandwidth consumption in the memory system when prefetches are not used or when they cause cache pollution. Figure 11 shows the bandwidth impact of prefetching in terms of *Bus Accesses per thousand Instructions (BPKI)*.¹⁰ Increasing the aggressiveness of the traditional stream prefetcher significantly increases the memory bandwidth consumption, especially for benchmarks where the prefetcher degrades performance. The proposed feedback directed prefetching mechanism reduces the aggressiveness of the prefetcher in these benchmarks. For example, in applu and ammp our feedback mechanism usually chooses the least aggressive prefetcher configuration and the least aggressive cache insertion policy as shown in Figures 7 and 9. This results in the large reduction in BPKI shown in Figure 11. Using the proposed feedback directed prefetching mechanism (*Dynamic Aggressiveness* and *Dynamic Insertion*) consumes 18.7% less memory bandwidth than the Very Aggressive traditional prefetcher configuration, while it provides 6.5% higher performance as shown in Section 5.3.

Table 6 shows the average performance and average bandwidth consumption of different traditional prefetcher configurations and the dynamic feedback directed prefetching. Compared to the traditional prefetcher configuration that consumes similar amount of memory bandwidth as feedback directed prefetching¹¹, the feedback directed prefetching mechanism provides 13.6% higher performance. Hence, incorporating our dynamic feedback mechanism into the stream prefetcher significantly increases the bandwidth-efficiency of the baseline stream prefetcher.

	No prefetching	Very Conservative	Middle-of-the-Road	Very Aggressive	Feedback-Directed
IPC	0.85	1.21	1.47	1.57	1.67
BPKI	8.56	9.34	10.60	13.38	10.88

Table 6. Average IPC and BPKI (Bus Accesses per thousand Instructions) for different prefetcher configurations.

¹⁰We use Bus Accesses (rather than the number of prefetches sent) as our bandwidth metric, because this metric includes the effect of L2 misses caused due to demand accesses as well as prefetches. If the prefetcher is polluting the cache, then the number of L2 misses due to demand accesses also increases. Hence, counting the number of bus accesses provides a more accurate measure of the memory bandwidth consumed by the prefetcher.

¹¹Middle-of-the-Road configuration consumes only 2.5% less memory bandwidth than feedback directed prefetching.

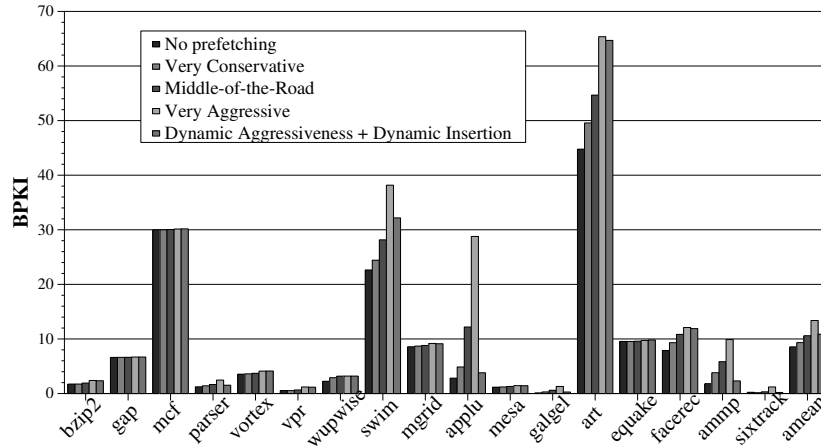


Figure 11. Effect of feedback directed prefetching on memory bandwidth consumption.

5.5. Hardware Cost and Complexity of Feedback Directed Prefetching

Table 7 shows the hardware cost of the proposed mechanism in terms of the required state. Feedback directed prefetching does not add significant combinational logic complexity to the processor. Combinational logic is required for the update of counters, update of the *pref-bits* in the L2 cache, update of the entries in the pollution filter, calculation of feedback metrics at the end of each sampling interval, determination of when a sampling interval ends, and insertion of prefetched blocks into appropriate locations in the LRU stack of an L2 cache set. None of the required logic is on the critical path of the processor. The storage overhead of our mechanism is less than 0.25% of the data-store size of the baseline 1MB L2 cache.

pref-bit for each tag-store entry in the L2 cache	16384 blocks * 1 bit/block = 16384 bits
Pollution Filter	4096 entries * 1 bit/entry = 4096 bits
16-bit counters used to estimate feedback metrics	11 counters * 16 bits/counter = 176 bits
pref-bit for each MSHR entry	128 entries * 1 bit/entry = 128 bits
Total hardware cost	20784 bits = 2.54 KB
Percentage area overhead compared to baseline 1MB L2 cache	2.5KB/1024KB = 0.24%

Table 7. Hardware cost of feedback directed prefetching.

5.6. Using only Prefetch Accuracy for Feedback

We use a comprehensive set of metrics –prefetch accuracy, timeliness, and pollution– in order to provide feedback to adjust the prefetcher aggressiveness. In order to assess the benefit of using timeliness as well as cache pollution, we evaluated a mechanism where we adapted the prefetcher aggressiveness based only on accuracy. In such a scheme, we increment the *Dynamic Configuration Counter* if the accuracy is high and decrement it if the accuracy is low. We found that, compared to this scheme that only uses accuracy to throttle the aggressiveness of a stream prefetcher, our comprehensive mechanism that also takes into account timeliness and cache pollution provides 3.4% higher performance and consumes 2.5% less bandwidth.

5.7. Feedback Directed Prefetching vs. Using a Prefetch Cache

Cache pollution caused by prefetches can be eliminated by bringing prefetched data into separate prefetch buffers [12, 10] rather than inserting prefetched data into the L2 cache. Figures 12 and 13 respectively show the performance and bandwidth consumption of the Very Aggressive prefetcher with different prefetch cache sizes - ranging from a 2KB fully-associate

prefetch cache to a 1MB 16-way prefetch cache.¹² The performance of the Very Aggressive prefetcher and the feedback directed prefetching mechanism when prefetched data is inserted into the L2 cache is also shown.

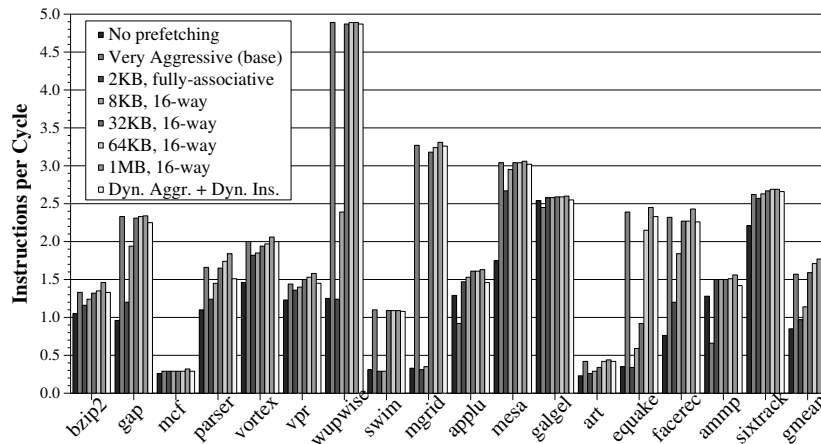


Figure 12. Performance comparison of using a prefetch cache vs. feedback directed prefetching.

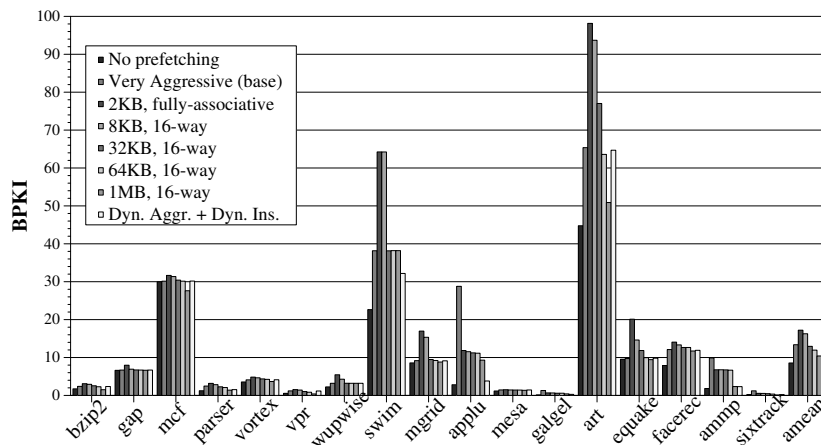


Figure 13. Bandwidth comparison of using a prefetch cache vs. feedback directed prefetching.

The results show that using small (2KB and 8KB) prefetch caches do not provide as high performance as inserting the prefetched data into the L2 cache. With an aggressive prefetcher and a small prefetch cache, the prefetched blocks are displaced by later prefetches before being used by the program - which results in lower performance with a small prefetch cache. However, larger prefetch caches (32KB and larger) improve performance compared to inserting prefetched data into the L2 cache because a larger prefetch cache reduces the pollution caused by prefetched data in the L2 cache while providing enough space for prefetched blocks.

Using feedback directed prefetching (both *Dynamic Aggressiveness* and *Dynamic Insertion*) that prefetches into the L2 cache provides 5.3% higher performance than that provided by augmenting the Very Aggressive traditional prefetcher configuration with a 32KB prefetch cache. The performance of feedback directed prefetching is also within 2% of the performance of the Very Aggressive configuration with a 64KB prefetch cache. Furthermore, the memory bandwidth consumption of feedback directed prefetching is 16% and 9% less than the Very Aggressive prefetcher configurations with respectively a

¹²In the configurations with a prefetch cache, a prefetched cache block is moved from the prefetch cache into the L2 cache if it is accessed by a demand load/store request. The block size of the prefetch cache and the L2 cache are the same and the prefetch cache is assumed to be accessed in parallel with the L2 cache without any adverse latency impact on L2 cache access time.

32KB and 64KB prefetch cache. Hence, feedback directed prefetching achieves the performance provided by a relatively large prefetch cache bandwidth-efficiently and without requiring as large hardware cost and complexity as that introduced by the addition of a prefetch cache that is larger than 32KB.

5.8. Effect of Feedback Directed Prefetching on a Global History Buffer Based C/DC Prefetcher

We have also implemented the proposed feedback directed prefetching mechanism on the C/DC (C-Zone Delta Correlation) variant of the Global History Buffer (GHB) prefetcher [9]. In order to vary the aggressiveness of this prefetcher dynamically, we vary the *Prefetch Degree*.¹³ Table 8 shows the aggressiveness configurations used for the GHB prefetcher. The dynamic feedback directed mechanism adjusts the configuration of the GHB prefetcher as described in Section 3.3.

<i>Dynamic Configuration Counter</i>	<i>Aggressiveness</i>	<i>Prefetch Degree</i>
1	Very Conservative	4
2	Conservative	8
3	Middle-of-the-Road	16
4	Aggressive	32
5	Very Aggressive	64

Table 8. Global History Buffer based C/DC prefetcher configurations.

Figure 14 shows the performance and bandwidth consumption of different GHB prefetcher configurations and the feedback directed GHB prefetcher using both *Dynamic Aggressiveness* and *Dynamic Insertion*. The feedback directed GHB prefetcher performs similarly to the best-performing traditional configuration (Very Aggressive configuration), while it consumes 20.8% less memory bandwidth. Compared to the traditional GHB prefetcher configuration that consumes similar amount of memory bandwidth as feedback directed prefetching (i.e. Middle-of-the-Road configuration), the feedback directed prefetching mechanism provides 9.9% higher performance. Hence, feedback directed prefetching significantly increases the bandwidth-efficiency of GHB-based delta correlation prefetching. Note that it is possible to improve the performance and bandwidth benefits of the proposed mechanism by tuning the thresholds used in feedback mechanisms to the behavior of the GHB-based prefetcher.

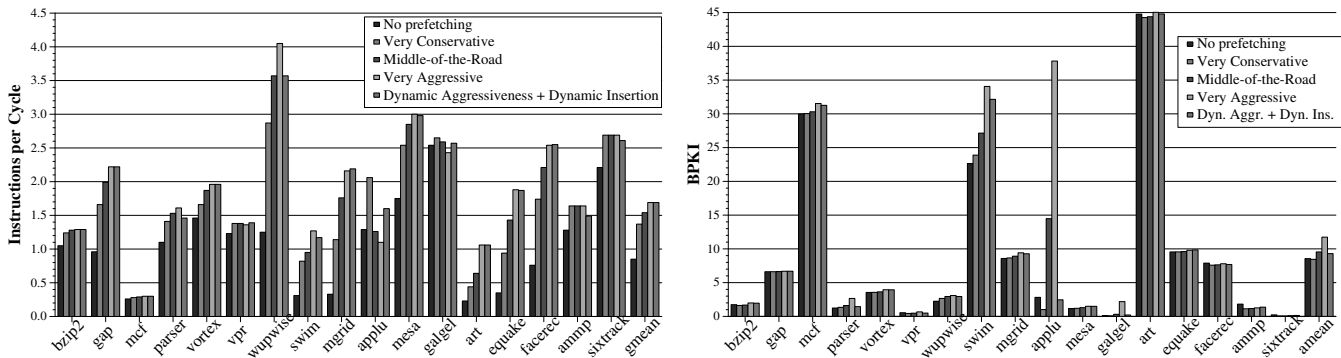


Figure 14. IPC performance (left) and memory bandwidth consumption in BPKI (right) of GHB-based C/DC prefetchers with and without using feedback directed prefetching.

¹³In the GHB-based prefetching mechanism, *Prefetch Distance* and *Prefetch Degree* are the same.

5.9. Effect of Feedback Directed Prefetching on a PC-Based Stride Prefetcher

We also evaluated the feedback directed prefetching mechanism on a PC-based stride prefetcher [1] and found that the results are similar to those achieved on both stream and GHB-based prefetchers. On average, using the feedback directed approach results in a 4% performance gain and a 24% reduction in memory bandwidth compared to the best-performing conventional configuration for a PC-based stride prefetcher. Due to space constraints, we do not present detailed graphs for these results.

5.10. Sensitivity to L2 Cache Size and Main Memory Latency

We also evaluated the sensitivity of the feedback-directed prefetching mechanism to different different cache sizes and memory latencies. In these experiments, we varied the L2 cache size keeping the memory latency at 500 cycles (baseline) and varied the memory latency keeping the cache size at 1MB (baseline). Table 9 shows the change in average IPC and BPKI provided by the feedback-directed prefetching mechanism over the best performing conventional prefetcher configuration. The feedback-directed prefetching mechanism provides better performance and consumes significantly less bandwidth than the best-performing conventional prefetcher configuration for all evaluated cache sizes and memory latencies. As memory latency increases, the IPC improvement of our mechanism also increases because increasing the effectiveness of the prefetcher becomes more important when memory becomes more of a performance bottleneck.

L2 Cache Size (memory latency = 500 cycles)						Memory Latency (L2 cache size = 1 MB)					
512 KB		1 MB		2 MB		250 cycles		500 cycles		1000 cycles	
Δ IPC	Δ BPKI	Δ IPC	Δ BPKI	Δ IPC	Δ BPKI	Δ IPC	Δ BPKI	Δ IPC	Δ BPKI	Δ IPC	Δ BPKI
0%	-13.9%	6.5%	-18.7%	6.3%	-29.6%	4.5%	-23.0%	6.5%	-18.7%	8.4%	-16.9%

Table 9. Change in IPC and BPKI with feedback-directed prefetching when L2 size and memory latency are varied.

5.11. Effect of Feedback Directed Prefetching on Other SPEC CPU2000 Benchmarks

Figure 15 shows the IPC and BPKI impact of feedback directed prefetching on the remaining 9 SPEC CPU2000 benchmarks that have less potential. We find that our feedback directed scheme provides 0.4% performance improvement over the best performing conventional prefetcher configuration (i.e. Middle-of-the-Road configuration) while reducing the bandwidth consumption by 0.2%. None of the benchmarks lose performance with feedback directed prefetching. Note that the best-performing conventional configuration for these 9 benchmarks is not the same as the best-performing conventional configuration for the 17 memory-intensive benchmarks (i.e. Very-Aggressive configuration). Also note that the remaining 9 benchmarks are not bandwidth-intensive except for fma3d and gcc. In gcc, the performance improvement of feedback directed prefetching is 3.0% over the Middle-of-the-Road configuration. The prefetcher pollutes the L2 cache and evicts many useful instruction blocks in gcc, resulting in very long-latency instruction cache misses that leave the processor idle. Using feedback directed prefetching reduces this negative effect by detecting the pollution caused by prefetch references and dynamically reducing the aggressiveness of the prefetcher.

6. Related Work

Even though mechanisms for data prefetching have been studied for a long time, dynamic mechanisms to adapt the aggressiveness of the prefetcher have not been studied as extensively as algorithms that decide what to prefetch. We briefly describe previous work in dynamic adaptation of prefetching policies in microprocessors.

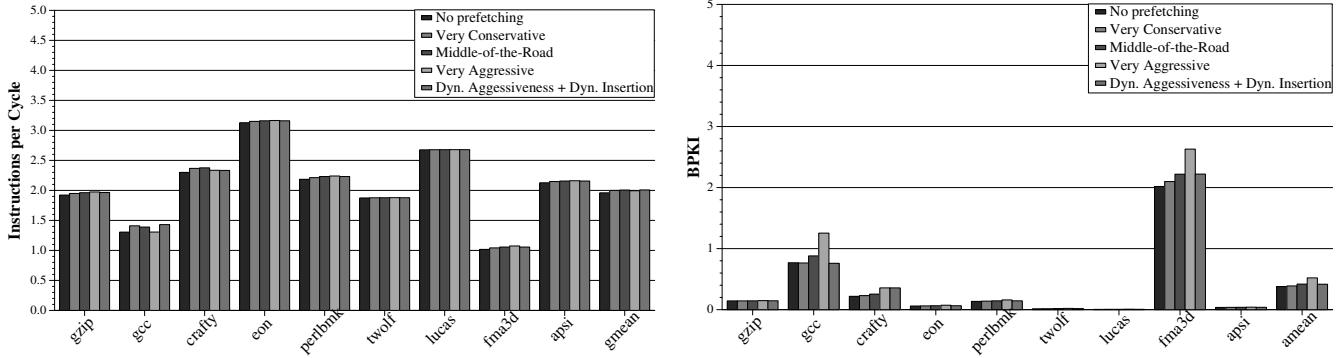


Figure 15. IPC performance (left) and memory bandwidth consumption in BPKI (right) of the stream prefetcher with and without using feedback directed prefetching.

6.1. Dynamic Adaptation of Data Prefetching Policies

The work most related to ours in adapting the prefetcher’s aggressiveness is Dahlgren et al.’s paper that proposed adaptive sequential (next-line) prefetching [4] for multiprocessors. This mechanism implemented two counters to count the number of sent prefetches (*counter-sent*) and the number of useful prefetches (*counter-used*). When counter-sent saturates, counter-used is compared to a static threshold to decide whether to increase or decrease the aggressiveness (i.e. *Prefetch Distance*) of the prefetcher. While Dahlgren et al.’s mechanism to calculate prefetcher accuracy is conceptually similar to ours, their approach considered only prefetch accuracy to dynamically adapt prefetch distance. Our proposal takes into account not only accuracy, but also prefetch timeliness and cache pollution to provide a comprehensive feedback mechanism. Also, Dahlgren et al.’s mechanism is designed for a simple sequential prefetching mechanism which prefetches up to 8 cache blocks following each cache miss. In this paper, we provide a generalized feedback-directed approach for dynamically adjusting the aggressiveness of state-of-the-art hardware data prefetchers. We show that our feedback mechanism improves the performance and bandwidth efficiency of the commonly used stream-based prefetchers as well as global-history-buffer based prefetchers and PC-based stride prefetchers.

Nesbit et al. [9] proposed tracking phases of a program execution. When the program enters a new phase of execution, the prefetcher is tuned based on the characteristics of the phase. This mechanism keeps track of phases by forming instruction working set signatures and associates a signature with the best-performing prefetcher configuration for that phase. In order to perform phase detection/prediction and identification of the best prefetcher configuration for a given phase, significant amount of extra hardware is needed. In comparison, our mechanism is simpler because it does not require phase detection or prediction mechanisms but instead adapts prefetcher behavior using three easy-to-compute metrics.

6.2. Cache Pollution Filtering

Charney and Puzak [3] proposed filtering L1 cache pollution caused by next-sequential prefetching and shadow directory prefetching from the L2 cache into the L1 cache. Their scheme associates a *confirmation bit* with each block in the L2 cache. The confirmation bit of a cache block is set if the block was used by a demand access when it was prefetched into the L1 cache the last time. When a prefetch request accesses the L2, it checks if the confirmation bit is set. If the confirmation bit is not set, the prefetch request is discarded, predicting that it will not be useful. For this scheme to be extended to prefetching from main memory to the L2 cache, a separate structure that maintains information about the blocks evicted from the L2 cache is required. This significantly increases the hardware cost of their mechanism. Our mechanism does not need to keep history information for evicted L2 cache blocks as it predicts how the prefetcher will perform based only on the three metrics

described.

Zhuang and Lee [23] proposed to filter prefetcher-generated cache pollution by using schemes similar to two-level branch predictors. Their mechanism tries to identify whether or not a prefetch will be useful based on past information about the usefulness of the prefetches generated to the same memory address or triggered by the same load instruction. This mechanism can reduce pollution only after recording the past behavior of a prefetched address or load instruction. In contrast, our mechanism does not require the collection of fine-grain information on each prefetch address or load address in order to vary the aggressiveness of the prefetcher.

Other approaches for cache pollution filtering include using a profiling mechanism to mark load instructions that can trigger hardware prefetches [21], and using compile-time techniques to mark dead cache locations so that prefetches can be inserted in dead locations [8]. In comparison to these two mechanisms, our mechanism does not require any software or ISA support and can adjust to dynamic program behavior even if it differs from the behavior of the compile-time profile. Lin et al. [14] proposed using density vectors to determine what to prefetch inside a region. This was especially useful in their model as they used very bandwidth-intensive scheduled region prefetching, which prefetches all the cache blocks in a memory region on a cache miss. This approach can be modified and combined with our proposal to further remove the pollution caused by blocks that are not used in a prefetch stream.

6.3. Cache Insertion Policy for Prefetches

Lin et al. [13] evaluated static policies to determine the placement in cache of prefetches generated by a scheduled region prefetcher. Their scheme placed prefetches in the LRU position of the LRU stack because they found that this policy resulted in good overall performance in their model. We found that, even though inserting prefetches in the LRU position reduces the cache pollution effects of prefetches on some benchmarks, it also reduces the positive benefits of aggressive stream prefetching on other benchmarks because useful prefetches—if placed in the LRU position—can be easily evicted from the cache in an aggressive prefetching scheme without providing any benefit. We have shown in Section 5.2 that statically placing prefetches in the LRU-4 position of the LRU stack provides good overall performance. Dynamically adjusting the insertion position of prefetched blocks based on the estimated pollution increases the average performance by 1.9% over the best static policy and by 18.8% over inserting prefetches in the LRU position.

7. Conclusion and Future Work

This paper proposed a feedback directed mechanism that dynamically adjusts the behavior of a hardware data prefetcher to improve performance and reduce memory bandwidth consumption. Our results show that using the proposed feedback directed mechanism increases average performance by 6.5% on the 17 memory-intensive benchmarks in the SPEC CPU2000 suite compared to the best-performing stream-based prefetcher configuration, while it consumes 18.7% less memory bandwidth. On the remaining 9 SPEC CPU2000 benchmarks, the proposed dynamic feedback mechanism performs as well as the best-performing conventional stream prefetcher configuration for those 9 benchmarks. Furthermore, the proposed feedback mechanism eliminates the performance loss incurred on some memory-intensive benchmarks due to prefetching. Over previous research in adaptive prefetching, our contributions are:

- We propose a comprehensive and low-cost feedback mechanism that takes into account prefetch accuracy, timeliness, and cache pollution caused by prefetch requests together to both throttle the aggressiveness of the prefetcher and to decide where in the cache to place the prefetched blocks. Previous approaches considered using only prefetch accuracy to determine the aggressiveness of simple sequential (next-line) prefetchers.

- We develop and evaluate a low-cost mechanism to estimate at run-time the cache pollution caused by hardware data prefetching.
- We propose and evaluate using comprehensive feedback mechanisms for state-of-the-art stream prefetchers that are commonly employed by today's high-performance processors. Our feedback-directed mechanism is applicable to any kind of hardware data prefetcher. We show that it works well with stream-based prefetchers, global-history-buffer based prefetchers and PC-based stride prefetchers. Previous adaptive mechanisms were applicable to only simple sequential prefetchers [4].

Future work can incorporate other important metrics, such as available memory bandwidth, estimates of the contention in the memory system, and prefetch coverage, into the dynamic feedback mechanism to provide further improvement in performance and further reduction in memory bandwidth consumption. Using genetic algorithms to determine the best decisions to make and the best thresholds to use in feedback directed prefetching can also increase the benefits provided by the proposed approach. The metrics defined and used in this paper could also be used as part of the selection mechanism in a hybrid prefetcher. Finally, even though this paper evaluated feedback mechanisms for data prefetchers, the mechanisms can be easily extended to instruction prefetchers.

References

- [1] J. Baer and T. Chen. An effective on-chip preloading scheme to reduce data access penalty. In *Proceedings of Supercomputing '91*, 1991.
- [2] B. H. Bloom. Space/time trade-offs in hash coding with allowable errors. *Communications of the ACM*, 13(7):422–426, 1970.
- [3] M. Charney and T. Puzak. Prefetching and memory system behavior of the SPEC95 benchmark suite. *IBM Journal of Research and Development*, 41(3):265–286, 1997.
- [4] F. Dahlgren, M. Dubois, and P. Stenström. Sequential hardware prefetching in shared-memory multiprocessors. *IEEE Transactions on Parallel and Distributed Systems*, 6(7):733–746, 1995.
- [5] J. D. Gindale. Buffer block prefetching method. *IBM Technical Disclosure Bulletin*, 20(2):696–697, July 1977.
- [6] G. Hinton, D. Sager, M. Upton, D. Boggs, D. Carmean, A. Kyker, and P. Roussel. The microarchitecture of the Pentium 4 processor. *Intel Technology Journal*, Feb. 2001. Q1 2001 Issue.
- [7] S. Iacobovici, L. Spracklen, S. Kadambi, Y. Chou, and S. G. Abraham. Effective stream-based and execution-based data prefetching. In *Proceedings of the 18th International Conference on Supercomputing*, 2004.
- [8] P. Jain, S. Devadas, and L. Rudolph. Controlling cache pollution in prefetching with software-assisted cache replacement. Technical Report CSG-462, Massachusetts Institute of Technology, 2001.
- [9] K. J. Nesbit, A. S. Dhodapkar, and J. E. Smith. AC/DC: An adaptive data cache prefetcher. In *Proceedings of the 13th International Conference on Parallel Architectures and Compilation Techniques*, pages 135–145, 2004.
- [10] N. P. Jouppi. Improving direct-mapped cache performance by the addition of a small fully-associative cache and prefetch buffers. In *Proceedings of the 17th International Symposium on Computer Architecture*, pages 364–373, 1990.
- [11] D. Kroft. Lockup-free instruction fetch/prefetch cache organization. In *Proceedings of the 8th Intl. Symposium on Computer Architecture*, 1981.
- [12] R. L. Lee, P.-C. Yew, and D. H. Lawrie. Data prefetching in shared memory multiprocessors. In *Proceedings of the Intl. Conference on Parallel Processing*, 1987.
- [13] W.-F. Lin, S. K. Reinhardt, and D. Burger. Reducing DRAM latencies with an integrated memory hierarchy design. In *Proceedings of the Seventh IEEE International Symposium on High Performance Computer Architecture*, pages 301–312, 2001.
- [14] W.-F. Lin, S. K. Reinhardt, D. Burger, and T. R. Puzak. Filtering superfluous prefetches using density vectors. In *International Conference on Computer Design*, pages 124–132. IEEE Computer Society, 2001.
- [15] O. Mutlu, H. Kim, D. N. Armstrong, and Y. N. Patt. An analysis of the performance impact of wrong-path memory references on out-of-order and runahead execution processors. *IEEE Transactions on Computers*, 54(12):1556–1571, Dec. 2005.
- [16] O. Mutlu, H. Kim, and Y. N. Patt. Techniques for efficient processing in runahead execution engines. In *Proceedings of the 32nd International Symposium on Computer Architecture*, 2005.
- [17] S. Palacharla and R. E. Kessler. Evaluating stream buffers as a secondary cache replacement. In *Proceedings of the 21st International Symposium on Computer Architecture*, pages 24–33, 1994.
- [18] J.-K. Peir, S.-C. Lai, S.-L. Lu, J. Stark, and K. Lai. Bloom filtering cache misses for accurate data speculation and prefetching. In *ICS '02: Proceedings of the 16th international conference on Supercomputing*, pages 189–198, New York, NY, USA, 2002. ACM Press.
- [19] A. J. Smith. Cache memories. *Computing Surveys*, 14(4):473–530, 1982.
- [20] L. Spracklen and S. G. Abraham. Chip multithreading: Opportunities and challenges. In *Proceedings of the Eleventh IEEE International Symposium on High Performance Computer Architecture*, pages 169–178, 2005.
- [21] V. Srinivasan, G. S. Tyson, and E. S. Davidson. A static filter for reducing prefetch traffic. Technical Report CSE-TR-400-99, University of Michigan Technical Report, 1999.
- [22] J. Tendler, S. Dodson, S. Fields, H. Le, and B. Sinharoy. POWER4 system microarchitecture. *IBM Technical White Paper*, Oct. 2001.
- [23] X. Zhuang and H.-H. S. Lee. A hardware-based cache pollution filtering mechanism for aggressive prefetches. In *Proceedings of the 32nd Intl. Conference on Parallel Processing*, pages 286–293, 2003.