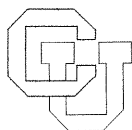# Hybridizing and Coalescing Load Value Predictors

**Martin Burtscher**
**Benjamin G. Zorn**

**CU-CS-903-00**

University of Colorado at Boulder

**DEPARTMENT OF COMPUTER SCIENCE**

# Hybridizing and Coalescing Load Value Predictors

Martin Burtscher
*Department of Computer Science*
*University of Colorado*
*Boulder, CO 80309-0430*
*burtsche@cs.colorado.edu*

Benjamin G. Zorn
*Microsoft Corporation*
*1 Microsoft Way*
*Redmond, WA 98052-6399*
*zorn@microsoft.com*

## Abstract

*Load value predictors usually require large amounts of state for retaining recently loaded values, which is particularly true for hybrid predictors. This paper presents several storage reduction techniques that significantly reduce the size of hybrids by sharing state between their components. Furthermore, our hybridization analysis shows that combining well-performing predictors does not always result in a good hybrid, whereas sometimes a relatively poor predictor can be a valuable addition to another predictor in a hybrid.*

*Detailed performance evaluations using a cycle-accurate microprocessor simulator running SPECint95 show that hybridizing can improve non-hybrids by close to forty percent over a wide range of sizes. With eleven kilobytes of state, our coalesced-hybrid yields a harmonic mean speedup of thirteen and fifteen percent with a re-fetch and a re-execute misprediction recovery mechanism, respectively, which is higher than the speedup of other predictors we evaluate, some of which are eight times larger.*

## 1. Introduction

Load instructions read data from memory rather than from the processor's fast register file. Because the memory hierarchy occasionally incurs long latencies, loads can take many cycles to execute, which slows down program execution. If the performance gap between CPUs and memory continues to widen, the load latency will become even longer. Unfortunately, load instructions are not only among the slowest but also among the most frequently executed instructions in current high-performance microprocessors. Hence, improving their execution speed can significantly boost the overall performance of a CPU.

Load instructions often fetch predictable sequences of values [11]. For instance, about half of all the load instructions in the SPECint95 benchmark suite retrieve the same value that they did the previous time they were exe-

cuted. Such behavior, which has been demonstrated explicitly on a number of architectures, is referred to as *value locality* [7, 11].

To exploit more of the existing load value locality, hybrid predictors have been proposed that combine several different predictors in one. A selector determines the best component for every prediction. Unfortunately, such hybrid predictors are often rather large [15, 22].

We devised two storage reduction optimizations that decrease the amount of state required by the well-performing *last n value* [4] and the *stride* predictor by a factor of two or more. We achieve this saving by letting the stride component reuse information already stored in the last *n* value component, making the former completely storage-less. In addition, we are able to shrink the last *n* value component by sharing 75% of the bits between the *n* values in each predictor line. Both techniques result in a significant decrease in predictor size yet have a negligible impact on the performance. The hybrid load value predictor we designed incorporates such a storage-less stride and a reduced-storage last three value predictor as well as a *register value* predictor [21], which is also storage-less. The resulting coalesced-hybrid is not only small but also highly effective. With only eleven kilobytes of state, it yields a speedup that surpasses the speedup of other, up to eight times larger predictors we considered both with a re-fetch and a re-execute misprediction recovery mechanism. Among predictors of similar size, the coalesced-hybrid outperforms other predictors by close to forty percent. Section 5.1 provides more results.

A detailed study of our hybrid's three main components (Section 5.2) reveals that they do indeed exploit distinct kinds of load value locality and thus contribute independently to the overall performance. This observation, which has also been made by Wang and Franklin [22] and others, indicates that predictors can be combined effectively to exploit a larger fraction of the existing load value locality. Building hybrid predictors may therefore be worthwhile in spite of their greater complexity. Our study further shows that not all predictors make good compo-

nents for a hybrid and, more surprisingly, that some predictors with a poor individual performance make a more valuable addition to a hybrid than other predictors with a good individual performance. Hence, detailed analyses are necessary to identify components that complement each other well.

The remainder of this paper is organized as follows: Section 2 introduces related work. Section 3 describes the storage reduction techniques and the architecture of our coalesced-hybrid load value predictor. Section 4 explains the evaluation methods. Section 5 presents the results. Section 6 concludes the paper with a summary.

## 2. Related Work

**Background**: To date, several categories of load value locality have been observed, including last value (sequences of identical values: e.g., 2, 2, 2, 2) [7, 11], stride (sequences of values with a constant offset between them: e.g., 1, 3, 5, 7, 9) [7, 16], last $n$ value (repetitions within the last $n$ values, e.g., 1, 2, 1, 2, 1, 2) [4, 10, 22], and finite context predictability (reoccurring arbitrary sequences of values: e.g., 1, 7, 3, ..., 1, 7, 3) [16]. Last value predictability is the simplest and most prominent kind of load value locality. Pure stride predictability (with a non-zero offset), on the other hand, is encountered only infrequently. Last $n$ value and finite context predictability have considerable potential but the latter is hard to exploit in small predictors. At least twenty percent of the dynamically executed load instructions cannot be predicted using any of the above schemes.

Like branch mispredictions, incorrect load value predictions necessitate a recovery process and thus incur a cycle penalty. Consequently, a load value predictor can actually slow down a processor instead of speeding it up if the percentage of incorrect predictions is large enough so that more cycles are added than saved. It is therefore important not to attempt a prediction if the prediction is likely to be incorrect. This is why almost all load value predictors are equipped with a *confidence estimator* (CE). Predictions are only allowed to take place if the estimated confidence that the prediction will be correct is high. There are two main approaches to confidence estimation in the current value prediction literature: saturating counters [11] and prediction outcome histories [2, 3, 5]. Both approaches have close counterparts in the branch prediction literature because confidence estimators are similar in design to branch predictors.

Saturating counters can count up and down within two boundaries, say zero and fifteen. If the counter has reached fifteen, counting up will not change its value. Likewise, counting down from zero leaves the counter at zero. The *bimodal* [12] confidence estimator uses such counters to record how many predictable values have

been seen in the recent past. The higher the value of the counter, the higher the confidence that the next load will be predictable since predictable load instructions do not frequently become unpredictable and vice-versa.

The *SAg* [23] confidence estimator represents an alternative approach. It works based on keeping a small history recording the most recent prediction outcomes (success or failure) [17]. Such histories consist of a short bit-pattern in which every bit indicates whether the corresponding prediction was correct. For instance, the left-most bit may represent the most recent prediction outcome, the next bit the second most recent outcome, etc. Every possible history pattern has a saturating counter associated with it to record the number of correct predictions that followed the corresponding history pattern in the recent past, thus assigning a confidence to each pattern.

Predicted load values allow the CPU to start processing the dependent instructions without having to wait for the memory access to complete, which improves the performance. *Speculative execution* enables the CPU to continue executing with a predicted value before the prediction outcome is known [18]. Because branch prediction requires a similar mechanism, most modern microprocessors already contain the necessary hardware to perform speculation.

Unfortunately, branch misprediction recovery hardware causes all the instructions that follow a misspeculated instruction to be purged and *re-fetched*. This is a very costly operation and makes a high prediction accuracy paramount. Unlike branches, which invalidate the entire execution path when mispredicted, mispredicted loads only invalidate the instructions that depend on the loaded value. In fact, even the dependent instructions per se are correct, they just need to be *re-executed* with the correct input value(s) [10]. A better recovery mechanism for load misspeculation therefore only re-executes the instructions that depend on the mispredicted load value. Such a recovery policy is less susceptible to mispredictions and favors a higher coverage, but may be hard to implement.

**Techniques**: Several research groups [4, 10, 22] have investigated last $n$ value predictability and noted its good performance. In this paper we show how the size of such a predictor can be reduced twofold by sharing the most significant bits among the $n$ values. We found that 75% of all the bits can be shared between the $n$ values in each predictor line basically without loss of performance.

Tullsen and Seng [21] present a register value predictor (Reg) that is storage-less except for its confidence estimator. It predicts that a load will fetch a value that is already in the target register of the load instruction before the load is executed. Since the predictor uses the CPU's register file as a source for values, it does not require any value storage in the predictor. This paper includes a per-

formance analysis of a register value predictor showing that it complements other predictors exceptionally well in a hybrid load value predictor. We further demonstrate how a stride predictor can also be made storage-less in combination with a last two value predictor.

**Other Predictors**: Lipasti et al. [11] designed a last value (LV) predictor with a bimodal confidence estimator. In prior work [2], we show that a SAg-based confidence estimator is able to improve their predictor's performance considerably. We therefore decided to also use SAg confidence estimators in our coalesced-hybrid predictor.

Sazeides and Smith [16] introduce the stride 2-delta (St2d) and the finite context method (FCM) predictor. The former maintains two separate strides and represents an improvement over the conventional stride predictor [7]. The stride used for making predictions is only updated if a new stride has been seen at least twice in a row, which reduces the number of mispredictions. A stride 2-delta predictor is included in our performance comparison in Section 5.1. Finite context method predictors retain short sequences of fetched load values. During a prediction they try to find the current sequence in their "database" and, if found, use the next value from the stored sequence to make a prediction.

**Hybrids**: Our performance comparison also includes a hybrid between a finite context method and a stride 2-delta predictor (St2d+FCM), as proposed by Rychlik et al. [15]. Their hybrid does not include any state reduction techniques.

Wang and Franklin designed a predictor that makes predictions based on the last four distinct values (LD4V) [22]. Their predictor uses a two-level access-pattern-based confidence estimator. In previous work, we show that it may not be necessary to store distinct values and propose a predictor that retains the last four values (L4V) independent of whether they are distinct or not [4]. In this paper we show that the size of such predictors can be reduced significantly by sharing the most of the bits between their sub-components.

Wang and Franklin further propose a hybrid predictor that combines their last four distinct value predictor with a stride predictor (LD4V+St). In the comparison section, we compare our predictor with both of Wang and Franklin's. In their hybrid, the stride component shares its base value with the last four distinct value component [22]. We now show that by storing not necessarily distinct values, the offsets required by the stride predictor can also be shared with the last $n$ value predictor, thus making the stride predictor completely storage-less.

## 3. Design of the Coalesced-Hybrid Predictor

Our predictor started out as a simple last value predictor with a SAg confidence estimator [2]. The first part of

Figure 3.1 (denoted as *Tag SAg LV*) shows an excerpt of four lines from such a predictor with eight-bit partial tags and ten-bit prediction outcome histories (the associated counters are not shown). It predicts that a load instruction will fetch the same value that it did the previous time it was executed, but the predicted value is only used if the partial tag matches and the confidence associated with the history in the selected predictor line is above a preset threshold.

In a previous publication [4], we show that even for moderate predictor sizes it is beneficial to reduce the height of a last value predictor in order to make it wider (e.g., a four times shorter predictor that stores the last four values in each line). Doing so increases the performance of the predictor without significantly changing the overall predictor size. The size does increase a little bit due to the replication of the second level of the SAg confidence estimator. The middle part of Figure 3.1 shows one line of a tagged SAg last four value predictor (*Tag SAg L4V*). The four sub-components operate independently and the component that reports the highest confidence is used for making a prediction. In case of a tie the component with the youngest value is selected [4]. Using the already present confidence information to guide the selection process eliminates the need for additional storage for selector related information [14, 15].
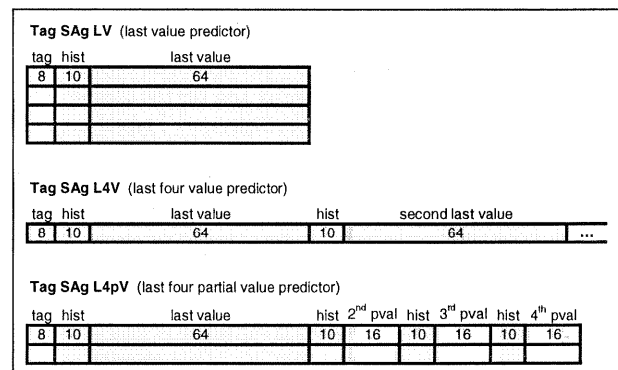


Figure 3.1: Architecture excerpts of three stages in the evolution of the coalesced-hybrid load value predictor. Only the tag and the first two of the *L4V*'s four components are shown.

We have already shown the last four value predictor to perform very well [4]. Now we improve this predictor further. All the improvements described in the remainder of this section are novel contributions of this paper.

First, we realized that the most significant bits of the four values within each predictor line are almost always identical. Hence, it suffices to store them only once instead of four times. Surprisingly, as many as 48 bits (or three quarters of all the bits) can be shared among the four values virtually without degrading the performance of the predictor while substantially reducing the storage re-

quirement. Our last four partial value predictor (*Tag SAg L4pV*) stores the full 64 bits of the most recently loaded (last) value but retains only the sixteen least significant bits of the three remaining values in each line. This reduces the predictor's size by about a factor of two. As a consequence, the predictor can now store twice as many values as its predecessor of the same size, which substantially improves the performance, in particular with small configurations. The last part of Figure 3.1 shows two lines of this predictor.

We then noticed that a last two value or wider predictor includes a "free" stride predictor. Stride predictors retain the last value and the difference (offset) between the last and the second to last value. The predicted value is the last value plus the offset. Our last four (partial) value predictor already retains the last value, and the stride can easily be computed on-the-fly out of the second to last value and the last value. The predicted value evaluates to the second last value subtracted from the last value multiplied by two (i.e., shifted to the right by one bit). The subtraction can be performed in parallel with the access to the second level of the confidence estimator since the two operations are independent. Except for the extra confidence estimator, the stride predictor is storage-less in combination with a last *n* value predictor (for $n \geq 2$).

Because the fourth component of the *L4pV* predictor hardly contributes to the overall performance (Section 5.5), we decided to leave it out, which provided a hardware saving that allowed us to add two additional confidence estimators, one of which we use for the storage-less stride predictor. We found Tullsen and Seng's register value predictor [21] to be an ideal candidate for the second confidence estimator since their predictor is also storage-less and only requires a confidence estimator.

We then added one more enhancement to our predictor that was first suggested by Bekerman et al. [1] and, independently, by Calder et al. [5]. They found that infrequently executed loads that alias with frequently executed loads evict useful predictor entries often enough to degrade the performance. According to their suggestion, we added one bit to the partial tags (which we termed *b-tags*) to indicate whether the last access to a given predictor line resulted in a tag miss. This bit makes it possible to prevent updating a predictor line after the first tag miss. Only allowing updates after at least two misses in a row very effectively prevents infrequently executed loads from being able to pollute the predictor.

Figure 3.2 shows the architecture of the coalesced-hybrid load value predictor with its storage-less register, storage-less stride, and reduced-storage last three partial value components (*St+Reg+L3pV*).

Every line of the predictor includes a nine-bit b-tag. The first of the five identical SAg confidence estimators (they each consist of an array of ten-bit histories "hist" and an array of three or four-bit saturating counters)

forms the register value component (*Reg*), which uses values from the CPU's register file for making predictions. The second confidence estimator belongs to the stride predictor (*St*) whose only other element is the adder since it uses values from the *L3pV* component. The remaining three confidence estimators, the 64-bit and the two 16-bit value fields form the last three partial value (*L3pV*) component. The predictor is pipelined over two stages similar to the way we pipelined the last four value predictor [4].
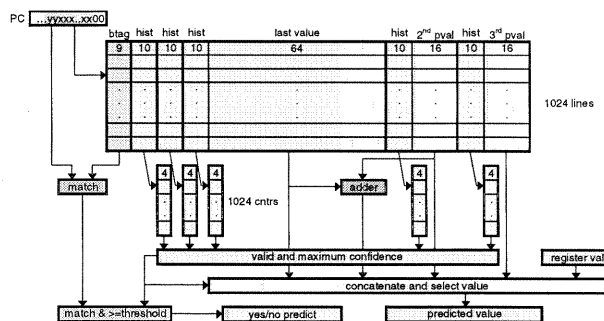


Figure 3.2: The architecture of our B-tag SAg *St+Reg+L3pV* "coalesced-hybrid" load value predictor.

The five sub-components operate independently and perform five value predictions and five confidence estimations in parallel. The value of the component reporting the highest confidence is used for making a prediction, but only if the confidence is above the preset threshold. To break ties, the five sub-components are prioritized from left to right, that is, the stride component has the highest priority, the register component has a medium priority, and the last three value component has the lowest priority. Within the last three value component, the more recent values have a higher priority [4]. Changing the prioritization order among the three main components has virtually no effect on the speedup. We use *St+Reg+L3pV* for no particular reason. The only minute performance difference we were able to detect is that prioritizing the *L3pV* component over the *Reg* component seems to be slightly disadvantageous.

When the predictor is updated, each component again makes a value prediction whose result is compared with the true load value. The confidence estimators are then updated based on the outcome of this comparison. At the same time, the values within the *L3pV* component are passed on to the next "older" sub-component and the true load value is copied into the 64-bit last value field.

# 4. Evaluation Methods

## 4.1 Benchmarks

We use the eight integer programs of the SPEC95 benchmark suite [19] with the provided reference input sets for our measurements. The executables were compiled using DEC GEM-CC with the highest optimization level "-migrate -O5 -ifo". The performed optimizations include common sub-expression elimination, split lifetime analysis, code scheduling, nop insertion, code motion and replication, loop unrolling, software pipelining, local and global inlining, inter-file optimization, etc. The binaries are statically linked to allow the linker to perform additional optimizations that reduce the number of runtime constants that are loaded during execution. These optimizations include most of the optimizations that OM [20] performs. The few floating point load instructions contained in the binaries are included in our measurements, loads to the zero-registers are ignored, and load immediate instructions are not taken into account since they do not access the memory and therefore do not need to be predicted. Table 4.1 gives relevant information about the eight benchmark programs.

Due to the detail of our simulations, each program is only executed for 300 million committed instructions on the simulator after having skipped over the initialization code in "fast-execution" mode. This fast-forwarding is important when only a section of a program's execution can be simulated because the initialization part of programs is usually not representative of the general program behavior [14]. The leftmost column of Table 4.1 shows the number of instructions that we skipped. *gcc* is simulated for 334 million committed instructions without skipping any instructions since this amounts to the full compilation of the *varasm* input-file.

| Information about the eight SPECint95 Benchmark Programs | | | | | |
|---|---|---|---|---|---|
| | million instrs | | percent | base | L1 load | L2 load |
| program | skipped | simul. | loads | IPC | miss-rate | miss-rate |
| compress | 6000 | 300 | 17.9% | 1.35 | 24.4% | 2.8% |
| gcc | 0 | 334 | 23.9% | 1.51 | 2.4% | 6.4% |
| go | 12000 | 300 | 24.1% | 1.44 | 1.4% | 15.3% |
| ijpeg | 1000 | 300 | 16.8% | 1.44 | 1.4% | 51.3% |
| li | 4000 | 300 | 25.5% | 1.99 | 5.4% | 0.6% |
| m88ksim | 1000 | 300 | 20.7% | 1.25 | 0.1% | 11.2% |
| perl | 1000 | 300 | 31.2% | 1.57 | 0.0% | 46.9% |
| vortex | 5000 | 300 | 23.6% | 2.89 | 2.2% | 10.2% |
| average | | | 22.9% | 1.68 | 4.7% | 18.1% |

Table 4.1: The first two columns show the number of instructions from the SPECint95 benchmark suite that are skipped and the number of instructions executed for our simulations, respectively. The third column lists the percentage of the executed instructions that are loads. The base IPC denotes the instructions per cycle that our baseline processor achieves on these programs. The last two columns give the L1 data-cache and the L2 cache load miss-rates.

In spite of the high optimization level and good register allocation, about every fifth instruction executed by these programs is a load. With an average IPC of 1.7, this amounts to one executed load instruction every 2.6 cycles.

With the exception of compress, the benchmark programs do not have very high L1 data-cache load miss-rates, making it harder for a load value predictor to be effective. On the other hand, the left half of Table 4.2 shows that only a relatively small number of load sites contribute most of the executed loads, implying that relatively small predictors should suffice to predict most of the executed loads.

The right half of Table 4.2 illustrates the load value predictability found in the eight benchmark programs. Register predictability "reg" indicates how often the target register of a load instruction already contains the value that the load is about to fetch. Last value predictability "lv" shows how often a load fetches a value that is identical to the previous value fetched by the same load instruction. Stride predictability "st2d" reflects how often a value is loaded that is identical to the last value plus the difference between the last and the second to last value fetched by the same load instruction. Last four value predictability "l4v" indicates how often a value is loaded that is identical to any one of the last four values fetched by the same load. Finally, finite context method predictability "fcm" shows how often a value is loaded that is identical to the value that followed the last time the same sequence of last four values was encountered (modulo a hash function). Note that, unlike reg, lv, st2d, and l4v, the fcm predictability results are implementation specific, i.e., they depend on the hash function.

| SPECint95 Quantile and Predictability Information | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | load sites that account for | | | | load value predictability (%) | | | | |
| program | Q100 | Q99 | Q90 | Q50 | reg | lv | st2d | l4v | fcm |
| compress | 62 | 56.5% | 45.2% | 14.5% | 9.0 | 40.4 | 65.8 | 41.3 | 35.9 |
| gcc | 34345 | 41.2% | 15.7% | 2.5% | 19.9 | 48.5 | 49.8 | 65.6 | 52.0 |
| go | 9619 | 40.2% | 17.9% | 2.7% | 9.2 | 45.9 | 47.2 | 64.0 | 44.6 |
| ijpeg | 2757 | 13.7% | 6.7% | 1.9% | 9.4 | 47.5 | 47.7 | 54.1 | 45.4 |
| li | 419 | 56.6% | 28.6% | 10.3% | 14.3 | 43.4 | 50.4 | 63.8 | 60.4 |
| m88ksim | 747 | 71.9% | 26.6% | 3.3% | 29.9 | 76.1 | 80.0 | 83.4 | 80.3 |
| perl | 1437 | 15.7% | 11.6% | 3.1% | 19.8 | 50.7 | 51.4 | 80.6 | 70.9 |
| vortex | 1973 | 48.6% | 18.0% | 2.8% | 17.8 | 65.7 | 65.3 | 78.6 | 66.9 |
| average | 6420 | 43.0% | 21.3% | 5.1% | 16.2 | 52.3 | 57.2 | 66.4 | 57.0 |

Table 4.2: The four quantile columns show the fraction of load sites that contribute the given percentage of executed loads. Q100 shows the number of loads that were executed at least once in absolute numbers. The remaining quantiles are relative to Q100. The five rightmost columns show the load value predictability in percent of executed loads found in each of the eight benchmark programs.

The predictability of the load instructions in all eight programs is quite high. On average, at least half of the executed load instructions are (theoretically) predictable using any method other than "reg".

## 4.2 Simulated Architecture

All our measurements are performed on the DEC Alpha AXP architecture [6] using the AINT simulator [13] with its cycle-accurate out-of-order back-end, which is configured to emulate a high-performance microprocessor similar to the DEC Alpha 21264 [9]. In particular, the simulated four-way superscalar CPU has a 128-entry instruction window, a 32-entry load/store buffer, four integer and two floating point units, a 64kB two-way set associative L1 instruction-cache, a 64kB two-way set associative L1 data-cache, a 4MB unified direct-mapped L2 cache, a 4096-entry BTB, and a 2048-line hybrid gshare-bimodal branch predictor. The modeled latencies are given in Table 4.3. The few operating system calls are executed but not simulated. Loads can only issue when all prior store addresses are known. The six functional units are fully pipelined and each unit can execute all operations in its class. Up to four load instructions are able to issue per cycle. This CPU represents our baseline.

| Instruction Type | Latency |
|---|---|
| integer multiply | 8-14 |
| conditional move | 2 |
| other int and logical | 1 |
| floating point multiply | 4 |
| floating point divide | 16 |
| other floating point | 4 |
| L1 load-to-use | 1 |
| L2 load-to-use | 12 |
| Memory load-to-use | 80 |

Table 4.3: The functional unit and memory access latencies (in cycles) used in our simulator.

To measure the speedup delivered by a load value predictor, the baseline CPU is augmented with the predictor in question and the resulting CPU performance is compared to the performance of the baseline processor. The pipelined load value predictions take place during the rename and issue-stage in the instruction pipeline and have a two-cycle latency. Note that even predicted loads perform a normal memory access. As soon as that access completes, the load value predictors are updated with the true load value. No speculative update is performed at the time of prediction. Out-of-order updates and updates from wrong-path loads are accurately modeled. Incorrect predictions may cause a conditional branch to transfer control to the wrong path.

Support for up to four predictor accesses per cycle (to match the issue-width of up to four loads per cycle) is provided by dividing all the predictors into four independent banks, as suggested by Gabbay and Mendelson [8]. Each bank can be thought of as an individual predictor one fourth the size. There is no communication between the banks, making it possible to operate them independently and in parallel. Since our simulator mimics a processor that fetches naturally aligned instructions, all

the load instructions that can possibly be fetched during the same cycle always go to distinct banks. To avoid conflicts between predictions and updates, updates are queued in a 16-entry FIFO queue (one per bank) and are dropped if the queue is full. The queue issues predictor updates at a rate of one per cycle whenever the queue is not empty and the corresponding predictor bank is idle.

## 5. Results

The following subsections describe the results. In Section 5.1 predictors from the literature are compared performance-wise with our coalesced-hybrid. Section 5.2 analyzes in detail the contributions of our hybrid predictor's components to the overall performance. In Section 5.3 our predictor is compared to oracles. Section 5.4 investigates the size of the partial values and Section 5.5 studies the width of the last $n$ value component. Due to space limitations, we only show average speedup results over the eight benchmarks and not the individual programs.

### 5.1 Comparison with Other Predictors

This section compares the harmonic-mean speedups over SPECint95 of several well-performing predictors from the literature and our own. The eight predictors we consider are: a tagged bimodal last value predictor (*Tag Bim LV*) [11], a tagged SAg last value predictor (*Tag SAg LV*) [2], a tagged SAg last four value predictor (*Tag SAg L4V*) [4], a tagged bimodal stride 2-delta predictor (*Tag Bim St2d*) [16], a tagged bimodal hybrid of a stride 2-delta and a finite context method predictor (*St2d+FCM*) [15], our coalesced-hybrid predictor (*St+Reg+L3pV*), a tagged last distinct four value predictor (*LD4V*) [22] with an access-pattern-based bimodal confidence estimator, and a hybrid between a *LD4V* and a stride predictor (*LD4V+Stride*) [22]. The performance of the individual components of our hybrid is discussed in the next section.

| Predictor Size in Kilobytes of State | | | |
|---|---|---|---|
| | 5-13kB | 19-27kB | 66-92kB |
| Tag Bim LV | 4.8 | 19.0 | 76.0 |
| Tag SAg LV | 5.8 | 21.1 | **82.6** |
| Tag SAg L4V | 7.3 | 21.5 | 78.5 |
| Tag Bim St2d | 5.8 | 23.0 | **92.0** |
| St2d+FCM | **11.4** | 23.2 | 65.8 |
| St+Reg+L3pV | 11.0 | 25.3 | 82.3 |
| LD4V | **12.6** | **26.6** | **82.3** |
| LD4V+Stride | **12.8** | **27.2** | **84.8** |

Table 5.1: The amount of state (in kilobytes) required by each of the eight predictors' three configurations. The bold numbers represent predictors that are larger than *St+Reg+L3pV*.

Since the predictors vary greatly in their architectures and complexities, they cannot be scaled to be of identical size. Consequently, we can only compare predictors of

similar sizes. In their base configurations, all eight predictors require between 19 and 27 kilobytes of state, which we believe is a realistic size for a first generation load value predictor. From these base-configurations we created two additional configurations for each predictor, a smaller one (by quartering the number of predictor lines) and a larger one (by quadrupling the number of predictor lines). The three size-ranges and the actual predictor sizes are shown in Table 5.1. The table gives the amount of state for a re-fetch architecture. For a re-execute architecture, some of the predictors require a little less state.

All the predictors are parametrizable in several dimensions and need to be configured to work well. To determine the setting that yields the highest speedup with our simulated CPU, we performed a detailed parameter space evaluation for most of the predictors. This evaluation includes varying the number of counter-bits, the prediction threshold, and the counter decrement (penalty) as well as the number of history-bits for the SAg-based predictors. For *LD4V* and *LD4V+Stride* we used the counter top and penalty given by Wang and Franklin [22], but we optimized the prediction threshold. *St2d+FCM* allows for many different ways of distributing the state over the two components. We performed a limited study to find a distribution that works quite well: 256 lines for the stride 2-delta component, 4096 lines (storing three nine-bit exclusive-or results representing the last three load values, which are shifted and xor'ed to compute the index into the second level) in the first level of the FCM component, and 1024 lines (that store the 64-bit load values) in the second level. Note that we use the counter parameters from *Tag Bim St2d* for *St2d+FCM*. Table 5.2 shows the base-configurations of the eight predictors.

no search for the optimal setting is performed. We use this approach to mimic what would happen if programs that are much larger or much smaller than the SPECint95 programs are run on these predictors. The intuition is that a larger program performs similarly on a load value predictor to a smaller program on a proportionately smaller version of the same predictor. Note that the number of lines in the first level of the FCM predictor is held constant when increasing and decreasing the size of the *St2d+FCM* predictor.

Figure 5.1 and Figure 5.2 present the harmonic-mean speedups of the eight predictors with a re-execute and a re-fetch misprediction recovery mechanism, respectively. Three speedup results are shown for each predictor corresponding to the three predictor sizes.
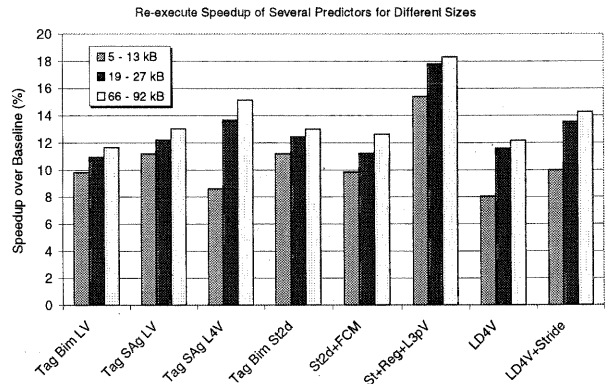


Figure 5.1: The re-execute speedup of several predictors for three size-ranges.

| | confid. estimator | value predictor | predictor lines | tag bits | hist bits | re-execute top | re-execute thr | re-execute pen | re-fetch top | re-fetch thr | re-fetch pen |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Tag Bim LV | bimodal | last value | 2048 | 8 | - | 7 | 3 | 2 | 15 | 9 | 7 |
| Tag SAg LV | SAg | last value | 2048 | 8 | 10 | 15 | 9 | 3 | 31 | 16 | 16 |
| Tag SAg L4V | SAg | last 4 values | 512 | 8 | 10 | 15 | 9 | 3 | 31 | 16 | 16 |
| Tag Bim St2d | bimodal | 2-delta stride | 2048 | 8 | - | 7 | 3 | 2 | 15 | 10 | 7 |
| St2d+FCM | bimodal | 2-delta str & finite cntxt | 256/4096 | 8 | - | 7 | 3 | 2 | 15 | 10 | 7 |
| St+Reg+L3pV | SAg | str & reg & last 3 pvals | 1024 | 8+1 | 10 | 15 | 8 | 4 | 31 | 16 | 16 |
| LD4V | bimodal | last distinct 4 values | 512 | 21 | - | 12 | 12 | 3 | 12 | 12 | 3 |
| LD4V+Stride | bimodal | last distinct 4 vals & str | 512 | 21 | - | 12/2 | 12/2 | 3/2 | 12/2 | 12/2 | 3/2 |

Table 5.2: The base-configurations of the eight predictors. Except in the coalesced-hybrid, strides are stored as eight-bit signed values. The coalesced-hybrid only stores the full 64 bits of the last value and the least significant 16 bits of the other two values. The counter top (*top*) represents the highest value that the saturating counters can reach. The lowest value is always zero. Predictions are made if the selected counter's value is at or above the given threshold (*thr*). If a component's prediction is correct, the corresponding counter is incremented by one and a one is shifted into the corresponding history, otherwise the counter is decremented by the given penalty (*pen*) and a zero is shifted into the history.

All the predictors are configured to work as well as possible in their base-configuration (19 to 27 kilobytes of state). Except for the number of predictor lines, the same parameters are used with the other two predictor sizes and
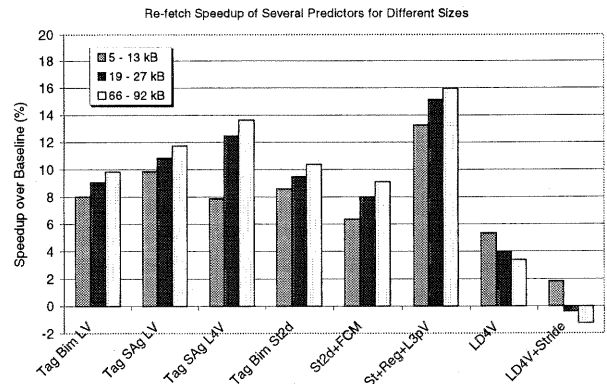


Figure 5.2: The re-fetch speedup of several predictors for three sizes ranges.

Our coalesced-hybrid predictor (*St+Reg+L3pV*) outperforms the other predictors both with a re-fetch and a re-execute misprediction recovery policy. More importantly, its re-fetch speedup exceeds the other predictors' re-execute speedup in all three size-ranges. Furthermore,

the performance of the smallest coalesced-hybrid configuration (requiring eleven kilobytes of state) surpasses the performance of the other predictors, including the ones from the largest size-range that require eight times more state. Only the last four value predictor, which is a predecessor of the coalesced-hybrid, is able to outperform the four times smaller coalesced-hybrid when re-fetch is used. This clearly shows that hybrid predictors do not necessarily have to be large to perform well and that coalescing the components in a hybrid predictor is a very effective technique to save state.

The good re-fetch speedup of our predictor is encouraging, in particular because it allows microprocessor designers to use the already existing branch misprediction hardware to recover from value mispredictions, which makes it less urgent to design and add a processor core that is capable of re-execution.

The performance of our predictor is certainly a result of a combination of factors, such as the b-tags, the inclusion of the register predictor, the SAg confidence estimators, and the incorporation of storage saving techniques. Nevertheless, the storage reduction is probably the key factor, as the following experiment illustrates. For comparison purposes, we built a hybrid between a stride 2-delta, a finite context method, and a register predictor with b-tags and SAg confidence estimators (*B-Tag SAg St2d+FCM+Reg*). Even with a global parameter optimization, this predictor is not able to reach the performance of the similarly sized coalesced hybrid.

Note that the performance of some of the predictors actually decreases with re-fetch when increasing the predictor size. Our investigation of this phenomenon revealed a somewhat surprising result. As it turns out, the smallest configuration of the two affected predictors both suffer significantly from aliasing. The confidence estimator detects this problem and inhibits the affected lines from making predictions. Consequently, the predictor only attempts relatively few predictions, which is reflected in its low performance when compared to the other predictors. The larger predictors suffer less from aliasing and the confidence estimator consequently allows more predictions to take place. Unfortunately, the CE also allows significantly more incorrect predictions, which more than offset the benefit of the additional correct predictions. Hence, the overall performance decreases as the predictors become larger.

Interestingly, the stride 2-delta predictor (*Tag Bim St2d*) performs better than the *St2d+FCM* hybrid. The reason is mostly that the small *FCM* component does not perform very well and therefore takes away valuable real-estate from the *St2d* component.

Among the predictors of a given size-range, the predictors with more components have fewer lines (i.e., are shorter) than the single-component predictors and are consequently more likely to experience capacity prob-

lems, in particular in the smallest configuration. The effect of the resulting aliasing can be observed in the two figures. The performance difference between the small and the middle configuration is significantly larger with the multi-component predictors (*Tag SAg L4V, LD4V,* and *LD4V+Stride*) than with the other predictors. Note that the coalesced-hybrid has more components than the *Tag SAg L4V* and the *LD4V* predictors, yet it is not affected as much by detrimental aliasing since the high degree of coalescing allows it to have twice the number of predictor lines (Section 3), which alleviates the capacity problem.

## 5.2 Component Contributions

To evaluate how much the individual components of the coalesced-hybrid contribute to the overall performance, we measured the speedup delivered by the hybrid's three main components in isolation, in pairs, and when all of them are used together. Figure 5.3 shows the results.
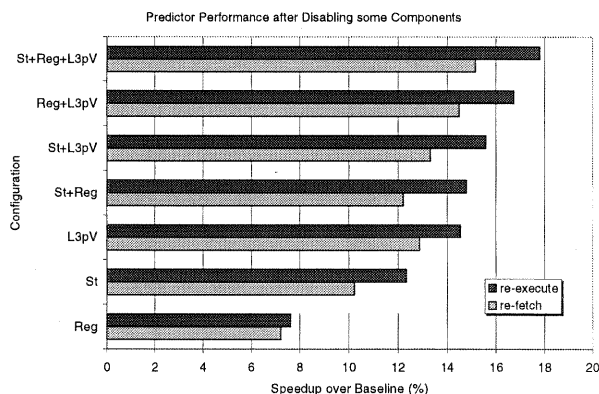


Figure 5.3: The performance of the coalesced-hybrid when disabling any combination of zero, one, or two of its three main components. Disabling means that the corresponding component cannot be selected for making a prediction and hence serves no purpose other than possibly providing shared state for some other component. Only the enabled components' names are shown.

As expected, using all three components yields the highest speedup both for re-fetch and re-execute. Using only one component results in the lowest speedup with one exception. The re-fetch performance of *L3pV* is better than the combined performance of *St+Reg*. Note, however, that the *L3pV* component is somewhat larger than the combined size of *St+Reg*.

Among the component pairs, *Reg+L3pV* performs best. Its performance is close enough to the performance delivered when all three components are combined that leaving the stride predictor out of the coalesced-hybrid may result in a more cost-effective implementation, in particular if the subtraction that is required by the stride component compromises the cycle time.

To better analyze the overlap between the individual speedup contributions, we used the seven configurations shown in Figure 5.3 as a set of seven equations and solved them for the distinct speedup contributions and overlaps. The result is depicted in the two Venn-diagrams in Figure 5.4. The displayed percentages add up to one hundred percent, which corresponds to the coalesced-hybrid's total harmonic mean speedup over the baseline processor (15.138% for re-fetch and 17.789% for re-execute).

The re-fetch Venn-diagram, for example, shows that the *St* component provides 4% of the total speedup that cannot be delivered by either one of the other two components. Similarly, *Reg* contributes 12% and *L3pV* 19% to the total speedup that none of the other components can provide. These are the three speedup contributions that are unique to the three components. The remaining four contributions are shared. The *St* and the *L3pV* component provide 29% of shared speedup, meaning that either one of the two components can provide this contribution, but the contribution does not increase if both components are used. *L3pV* and *Reg* add 1% of shared speedup. The shared speedup between *St* and *Reg* is -1%, indicating that the two components interfere negatively and occasionally prohibit one another from making a correct prediction. The contribution that is shared among all three components is 36%.
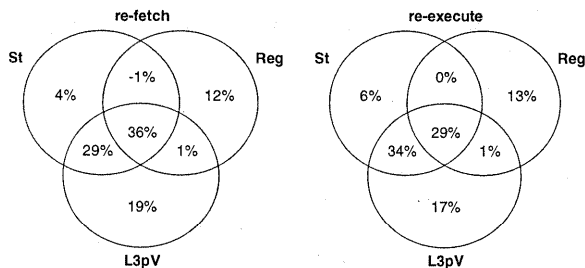


Figure 5.4: Venn-diagrams showing the overlap and the individual speedup contributions of the three main predictor components.

Approximately a third of the speedup can be delivered by any one of the three components and another third by either the stride or the last three partial value component. More importantly, 17% to 19% of the speedup can only be provided by the *L3pV* component and 12% to 13% only by the register component. The stride component, on the other hand, only delivers 4% to 6% of speedup that cannot be attained by either one of the other two components. Moreover, the intersection of the register predictor with the other two components is relatively small, indicating that the *Reg* component is able to predict a rather distinct set of loads. This observation is consistent with the results from Figure 5.3, which show that the register component by itself does not perform very well but makes

a strong combination with the *L3pV* component exactly because it can handle important loads that the *L3pV* predictor cannot. Note that we did not use profiling to change the register allocation, which can significantly improve the performance of the *Reg* predictor [21], yet we already get a substantial benefit from including a register value predictor.

It is surprising that the *Reg* predictor, which performs poorly when used by itself, complements the *L3pV* component much better than the *St* predictor with its decent individual performance. This nicely illustrates the importance of detailed component analyses to find cooperative components for building hybrids and that unconventional predictors with a poor individual performance can make a valuable addition to a hybrid predictor.

## 5.3 Comparison with Oracles

In this section we compare the coalesced-hybrid with versions of itself that contain oracles to demonstrate how much of the existing potential the predictor can reap.

The first predictor (*normal*) in Figure 5.5 represents the coalesced-hybrid in its conventional and implementable form as described in Section 3. It does not include an oracle. The first oracle (*perf-ce/sel*) represents the same predictor except it incorporates a perfect confidence estimator and a perfect selector. This means that, whenever possible, the component that will make a correct prediction is selected and forced to make a prediction. If no such component exists, no prediction is attempted. Therefore, the oracle never makes a misprediction. The second oracle (*perf-all*) simply predicts every executed load with the correct value. There are no mispredictions and, as opposed to *perf-ce/sel*, this oracle never decides not to make a prediction. Figure 5.5 shows the speedups delivered by the oracle-less predictor and the two oracles.
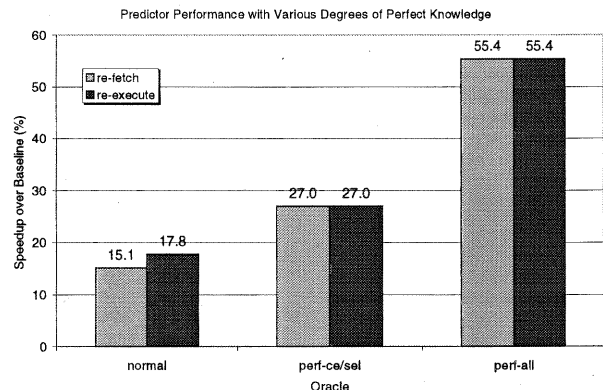


Figure 5.5: Re-fetch and re-execute speedups of the coalesced-hybrid with different degrees of perfect knowledge.

A perfect confidence estimator in combination with a perfect selector (*perf-ce/sel*) results in a significant increase in speedup over the conventional predictor (*normal*). A more detailed analysis revealed that both the selection mechanism and the confidence estimator of the normal predictor are far from perfect. In particular, the coalesced-hybrid's imperfect CE is rather conservative and inhibits a considerable number of predictions that would be correct. Since the CE setting we use is the result of a global optimization and yields one of the highest possible speedups, trading off missing potentially correct predictions for reducing the number of incorrect predictions must be advantageous in the CPU we are simulating. Overall, the coalesced-hybrid's confidence estimator and selector (*normal*) are able to reap 56% to 63% of the theoretically possible speedup (*perf-ce/sel*) for this predictor.

A comparison with the perfect load value predictor (*perf-all*), however, shows that there is still significant potential for improvement left. Our predictor only yields 27% to 31% of the speedup that can theoretically be attained with load value prediction (*perf-all*). Comparing *perf-all* with *perf-ce/sel* shows that the coalesced-hybrid only contains the necessary information to reach about half the possible speedup. This large gap suggests that there exists significant opportunity for new and different prediction methods to improve the performance beyond that of existing methods. It is, however, unclear how much of the remaining potential can be realized because the fraction of unpredictable loads is unknown.

## 5.4 Partial Value Size

To determine how many bits can be shared among the sub-components of the last three partial value predictor without overly impacting the performance, we present Figure 5.6. It shows the speedup of the coalesced-hybrid with varying number of bits in the partial value fields.
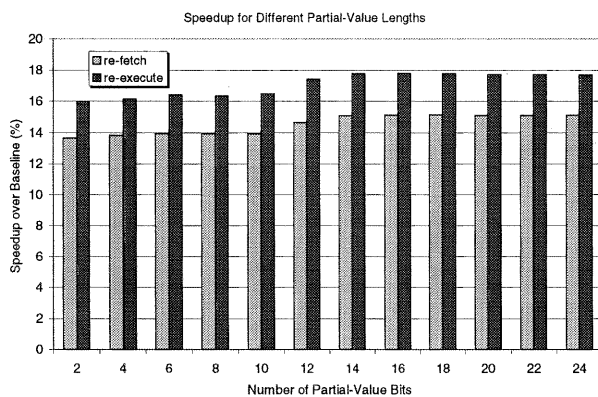


Figure 5.6: The speedup of the coalesced-hybrid when the number of bits in the partial values is varied.

The predictor performs well even with extremely short partial values. Unfortunately, this is mostly because the last value and the register components are not affected by shortening the partial values and hence their performance is independent of the size of the partial values.

Nevertheless, a performance jump can be observed between ten and fourteen-bit partial values. It appears that fourteen-bit values can capture a substantially larger fraction of the occurring values in the SPECint95 programs than ten-bit values. At the same time it looks like fourteen bits are enough to handle the vast majority of (predictable) values since there is almost no speedup improvement when further increasing the number of bits in the partial values. Likewise, the performance does not degrade much below ten bits. We use sixteen bits to be safely on the upper plateau.

We also tried adding a valid bit to the two partial value fields to indicate whether the concatenation of the last value's 48 most significant bits with the 16-bit partial value yields the correct 64-bit value, thus allowing only "valid" components to make a prediction. As it turns out, the confidence estimator already keeps track of this information and the valid bits are superfluous.

We further tried using 16-bit signed offsets that had to be added to the shared 48 most significant bits. However, this increased the complexity of the predictor and resulted in poorer performance than using 16-bit unsigned values that only have to be concatenated with the shared most significant bits.

## 5.5 Predictor Width

To determine how wide the last $n$ partial value component should be, we present Figure 5.7. It shows the performance of the coalesced hybrid with differently sized last $n$ partial value components.
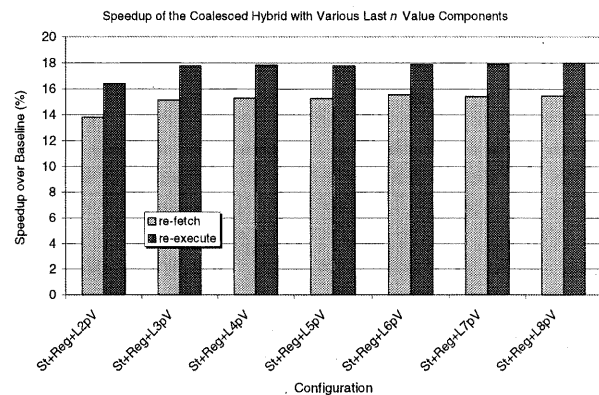


Figure 5.7: The speedup of the coalesced-hybrid when the width of the last $n$ partial 16-bit value component is varied.

Clearly, retaining three to four last values per line is

sufficient to reap almost all the potential. This result is particularly surprising because the predictors used in the figure are not scaled to the same size but become larger as the width increases. Consequently, we use a last three partial value component to keep the predictor's size and the number of components small while still reaping most of the performance.

Note that this result indicates that hybridization is more important than making the predictor's components wider to improve their individual performance [4].

The slight performance fluctuations (e.g., *L7pV*'s re-fetch performance is marginally lower than *L6pV*'s) are due to increased negative interference between the predictor's components.

## 6. Summary and Conclusions

In this paper we perform a detailed performance evaluation of the components of a hybrid load value predictor and describe two powerful state reduction techniques that allowed us to design a very effective hybrid that requires only a small amount of state.

Our study of three hybrid components (a stride, a register value, and a last three value predictor) shows that different components exploit different kinds of load value locality and that they contribute independently to the overall performance.

Our study further shows that care must be taken when selecting predictors as components for a hybrid because some predictors with a poor individual performance make a more valuable addition to a hybrid than other predictors with a good individual performance. To identify components that complement each other well, performance analyses are most likely unavoidable.

To reduce the often large storage requirement of hybrid predictors, we devised two storage reduction techniques that decrease the amount of state required by a last $n$ value and a stride predictor by a factor of two or more. We achieve this saving by having the last $n$ value component provide all the information that the stride component needs, making the latter storage-less. In addition to that, the size of the last $n$ value predictor is reduced by sharing most of the bits among the $n$ values in each predictor line. Both techniques result in a substantial decrease in predictor size virtually without impacting the performance.

The hybrid load value predictor we designed incorporates such a storage-less stride and a reduced-storage last three value predictor as well as a storage-less register value predictor. Cycle-accurate pipeline-level simulations of a four-way superscalar out-of-order CPU with many different load value predictors show that our predictor outperforms other predictors by almost forty percent over a large range of sizes. In the smallest configuration we investigated, which requires eleven kilobytes of state, our coalesced-hybrid yields a speedup with a re-fetch and with a re-execute misprediction recovery mechanism that surpasses the speedup of other predictors from the literature, some of which are eight times larger.

We believe that a large fraction of the value locality found in short load value sequences has been captured. However, it remains an open research question whether longer sequences contain significant additional predictability and how much of it can be extracted efficiently using existing and new prediction techniques. In future work we will investigate how approaches like the finite context method can be incorporated into our predictor to exploit even more of the existing load value locality.

## Acknowledgments

## References

[1]  M. Bekerman, S. Jourdan, R. Ronen, G. Kirshenboim, L. Rappoport, A. Yoaz, U. Weiser. "Correlated Load-Address Predictors". *26th International Symposium on Computer Architecture*. May 1999.

[2]  M. Burtscher, B. G. Zorn. *Load Value Prediction Using Prediction Outcome Histories*. Unpublished Technical Report CU-CS-873-98, University of Colorado at Boulder. October 1998.

[3]  M. Burtscher, B. G. Zorn. "Prediction Outcome History-based Confidence Estimation for Load Value Prediction". *Journal of Instruction-Level Parallelism*, Vol. 1. May 1999.

[4]  M. Burtscher, B. G. Zorn. "Exploring Last $n$ Value Prediction". *1999 International Conference on Parallel Architectures and Compilation Techniques*. October 1999.

[5]  B. Calder, G. Reinmann, D. M. Tullsen. "Selective Value Prediction". *26th International Symposium on Computer Architecture*. May 1999.

[6]  Digital Equipment Corporation. *Alpha Architecture Handbook*. 1992.

[7]  F. Gabbay. *Speculative Execution Based on Value Prediction*. EE Department Technical Report #1080, Technion - Israel Institute of Technology. November 1996.

[8]  F. Gabbay, A. Mendelson. "The Effect of Instruction Fetch Bandwidth on Value Prediction". *25th International Symposium on Computer Architecture*. June 1998.

[9]  R. E. Kessler, E. J. McLellan, D. A. Webb. "The Alpha 21264 Microprocessor Architecture". *1998 International Conference on Computer Design*. October 1998.

[10] M. H. Lipasti, J. P. Shen. "Exceeding the Dataflow Limit via Value Prediction". *29th International Symposium on Microarchitecture*. December 1996.

[11] M. H. Lipasti, C. B. Wilkerson, J. P. Shen. "Value Locality and Load Value Prediction". *Seventh International Conference on Architectural Support for Programming Languages and Operating Systems*. October 1996.

[12] S. McFarling. *Combining Branch Predictors*. WRL Technical Note TN-36, Digital Western Research Laboratory, Palo Alto. June 1993.

[13] A. Paithankar. *AINT: A Tool for Simulation of Shared-Memory Multiprocessors*. Master's Thesis, University of Colorado at Boulder. 1996.

[14] G. Reinman, B. Calder. "Predictive Techniques for Aggressive Load Speculation". *31st International Symposium on Microarchitecture*. December 1998.

[15] B. Rychlik, J. Faistl, B. Krug, J. P. Shen. "Efficacy and Performance Impact of Value Prediction". *1998 International Conference on Parallel Architectures and Compilation Techniques*. October 1998.

[16] Y. Sazeides, J. E. Smith. "The Predictability of Data Values". *30th International Symposium on Microarchitecture*. December 1997.

[17] E. Sprangle, R. Chappell, M. Alsup, Y. Patt. "The Agree Predictor: A Mechanism for Reducing Negative Branch History Interference". *24th International Symposium on Computer Architecture*. June 1997.

[18] J. E. Smith, G. S. Sohi. "The Microarchitecture of Superscalar Processors". *Proceedings of the IEEE*. 1995.

[19] *SPEC CPU'95*. August 1995.

[20] A. Srivastava, D. Wall. "A Practical System for Intermodule Code Optimization at Linktime". *Journal of Programming Languages* 1(1). March 1993.

[21] D. Tullsen, J. Seng. "Storageless Value Prediction Using Prior Register Values". *26th International Symposium on Computer Architecture*. May 1999.

[22] K. Wang, M. Franklin. "Highly Accurate Data Value Prediction using Hybrid Predictors". *30th International Symposium on Microarchitecture*. December 1997.

[23] T. Y. Yeh, Y. N. Patt. "A Comparison of Dynamic Branch Predictors that use Two Levels of Branch History". *20th International Symposium on Computer Architecture*. May 1993.