

Address-Free Memory Access Based on Program Syntax Correlation of Loads and Stores

Lu Peng, Jih-Kwon Peir, Qianrong Ma, and Konrad Lai

Abstract—An increasing cache latency in next-generation processors incurs profound performance impacts in spite of advanced out-of-order execution techniques. One way to circumvent this cache latency problem is to predict load values at the onset of pipeline execution by exploiting either the load value locality or the address correlation of stores and loads. In this paper, we describe a new load value speculation mechanism based on the program syntax correlation of stores and loads. We establish a *Symbolic Cache (SC)*, which is accessed in early pipeline stages to achieve a zero-cycle load. Instead of using memory addresses, the SC is accessed by the encoding bits of base register ID plus the displacement directly from the instruction code. Performance evaluations using SPEC95 and SPEC2000 integer programs on SimpleScalar simulation tools show that the SC achieves higher prediction accuracy in comparison with other load value speculation methods, especially when hardware resources are limited.

I. INTRODUCTION

Today’s high-performance processor pipeline permits overlapping instruction execution to achieve more than one Instruction Per Cycle (IPC) average execution rate. The available Instruction-Level Parallelism (ILP) constrains this parallel execution because dependent instructions must wait for the data produced by the source instructions. The severity, in terms of execution delays, depends primarily on the speed that the producer instruction can generate the needed data.

Memory load latency presents a classical pipeline bottleneck even when the data is located in the first-level cache (L_1). Usually, the load data from L_1 is not ready until late stages of the pipeline while the dependent instruction requires the data at an earlier stage. This load-to-use delay exacerbates in recent high-performance microprocessors in which multi-cycle, first-level caches become the norm [21], [24], [23], [14], [12]. As the cache size, clock frequency, and complexity of microarchitecture continue to increase in next-generation processors, it is estimated that the L_1 cache accesses may consume two to five cycles [2]. This increasing load latency from caches will further lengthen the load-to-use delay and will have profound performance impacts in spite of advanced out-of-order execution techniques [2], [3], [18]. Simulations using SPEC2000 integer benchmarks running on the out-of-order SimpleScalar model [4] have shown that each cycle reduction of the L_1 access delay improves the IPC by 5–10% [18].

Lu Peng and Jih-Kwon Peir are with Department of Computer and Information Science and Engineering, University of Florida, Gainesville, FL 32611 Email: {lpeng,peir}@cise.ufl.edu

Qianrong Ma is with Server Technology, Oracle Corporation, 200 Oracle Pkwy 20p6, Redwood City, CA 94065 Email: qianrong.ma@oracle.com

Konrad Lai is with Microprocessor Research, Intel Lab, Intel Corporation, 5200 NE Elam Young Parkway, Hillsboro, OR 97124 Email: konrad.lai@intel.com

In Figure 1, a conceptual out-of-order execution pipeline is partitioned into two phases. First, an instruction is fetched, decoded, renamed, and issued through the *front-end* of pipeline stages. Afterwards, the register operands are read and the instruction is executed (including memory access) and committed through the *back-end* of pipeline stages. In order to be stall-free, a source instruction must produce the data before its dependent executions. In other words, a critical producer, when it is fetched and issued at the same cycle as its dependent instructions, needs to generate the result in the front-end of the pipeline to avoid any stall of its dependents. Such a dependent stall-free memory load instruction is called a *zero-cycle* load.

Fig. 1. Processor pipeline and dependent stall-free point

There have been several attempts to achieve a zero-cycle load by predicting and speculating the load value [15], [16], [26], [22], [25], [11], [5] or the load address [9], [10], [6], [2] in the front-end of the processor pipeline. Both load value and load address predictions generally suffer a low prediction accuracy. For address predictions, a lengthy cache access is still required that may delay the load dependents even if the predicted load address is correct.

In this paper, we exploit a new avenue to speculatively obtain the load value in front-end stages of the pipeline. First, we observe that store-load and load-load correlations are established in software and often displayed in the program syntax in the form of a base register ID plus a displacement value. Therefore, it is reasonable to use part of the store/load encoding bits (base register ID + displacement) directly to capture such correlations. Second, applications exhibit spatial locality among memory references. Such locality can also be observed in the program syntax when nearby loads or stores differ only by a small displacement value. Therefore, it is beneficial to establish store/load dependences on a large block granularity to capture the spatial data reference locality.

The syntax correlation holds when the content of the base register remains unchanged. This property exists in various program constructs such as accessing global and local variables, saving/restoring registers during procedure/function calls, referencing different records using the same pointer in linked data structures, accessing array elements in loop iterations with/without loop unrolling, etc. We also observe that the base address may stay the same even when the base register is updated between two memory references. This is due to a lack of sufficient registers, an uncertainty of future execution paths, or a traversal through different procedures that requires a base register to be saved and restored before the next usage.

Based on these observations, we propose a *Symbolic Cache (SC)*. An SC is a small separate data cache that is accessed at early front-end stages using certain encoding bits directly from memory instructions. The speculative data retrieved from the SC can trigger the execution of dependent instructions to avoid any delays. Performance evaluations using SimpleScalar tools and SPEC95/SPEC2000 integer programs show that the average prediction accuracy reaches over 70% using small SCs. This accuracy is generally higher than other data speculation methods, especially when hardware resources are limited for constructing extra caches and tables. The remaining paper is organized as follows. A few related work on hiding cache latency will be given in the next section. The motivations and important observations for the proposed method will be described in Section 3. This is followed by discussions of design and related issues for establishing the SC in Section 4. In Section 5, performance evaluations of three data speculation methods are given. Several design parameters for the SC are also evaluated. Finally, Section 6 concludes the paper.

II. RELATED WORK

The most aggressive load value speculation is to predict the value at the onset of pipeline execution. A load-value history table is established and accessed using the Program Counter (PC) of the load. This scheme allows loads bypassing caches completely to achieve a zero-cycle load. A value prediction can be successful if the value is repeated from the previous execution of the load [15], [16], [26], or the load value is followed certain recurrence patterns [22]. However, the lack of a close correlation between the instruction address and the value of the load makes it difficult to achieve a high prediction accuracy [15], [16], [26], [22], [11], [5].

Another way to circumvent pipeline hazards caused by the cache latency is to predict the load address at the onset of pipeline execution so that a cache access can start speculatively without going through the normal decode, rename, and address generation stages [9], [10], [6], [2]. Existing address prediction methods exploit regular patterns such as stride-based address patterns, and irregular but repeated patterns such as addresses for traversing link-based data structure. However, the difficulty remains of predicting a significant portion (over 30% [2]) of load addresses that do not fall into these two categories. In a recent proposal, dynamic dependence links were established between the instruction which updates a register to the instruction where the register is used as the base register [8]. Once the updated value is available, the dependent load address can be calculated early and more accurately. However, the lengthy cache access is required still, even with a correct address.

Memory renaming techniques establish dynamic dependence correlations between stores and loads [25]. A separate storage element called a *value file (VF)* is used to save the correlated data. When a memory load instruction is fetched, an indirect access to the value file based on the PC of the load can retrieve the data without going through a lengthy cache access. Studies show that there are many more loads that consume the value from the same producer than those loads which repeat the same value or address from the previous instance of the same load. Therefore, there is a better chance to obtain the correct

load value by using memory renaming through the VF rather than based on the load value/address locality. This approach, however, requires additional hardware to establish the correct dependence links among stores and loads. The load value cannot be accurately predicted before such a correlation has established dynamically. A similar idea has been exploited to dynamically establish store-load [19] and load-load [20] *associations*. A small *synonym file* which keeps the correlated data can be indirectly accessed by the PC of the load.

Recently, another early load address resolution technique for deep-pipelined machines has been proposed [3]. The authors observed that the addresses for certain types of memory loads, such as stack access, constant, or stride-based memory access, have regular increment/decrement patterns. By tracking the registers used for this type of load, register updates can be computed at the decode stage. As a consequence, the dependent load can start the address generation and cache access earlier after the load is decoded. Although non-speculative, this approach is limited to memory loads with certain address patterns. Also, the lengthy cache access is still required.

There have been other attempts to achieve fast cache accesses. The real cache index bit prediction based on the base register content enables parallel address translation and cache access [13]. Due to small offset values, the zero-cycle load technique [1] uses a simple carry-free adder for fast approximation of the load address. To avoid speculative address calculations, a special compiler-directed register is added in [7] to save the content of the base register for the next load so that the load address can be calculated in the decode stage. The SAM cache [17] uses the base address and the offset separately to access the cache directly. Although all these techniques achieve fast cache access, their impact in hiding the long cache latency on deep-pipelined microarchitectures is rather limited.

The proposed *Symbolic Cache (SC)* has several advantages over existing cache latency hiding methods. First, the SC can handle any type of loads, address patterns, or special usages of base registers. Second, unlike address predictions or register tracking, loads through the SC can bypass the address generation and cache access completely to achieve a zero-cycle load. This is similar to the value prediction method. However, instead of being based upon the history of the load values, the SC captures store/load syntax correlations with higher accuracy. Third, unlike the memory renaming technique, where the store/load correlation is established dynamically by the hardware, the store/load correlation is directly obtained from the instruction encoding bits to simplify the hardware requirement. In addition, the SC can capture spatial locality among memory references.

III. SYNTAX CORRELATION OF MEMORY REFERENCES

The foundation of the *Symbolic Cache (SC)* is based on store-load and load-load correlations from the program syntax in the form of a base register ID and a displacement value. This simple memory reference syntax also exhibits spatial locality. In this section, we will provide two programming examples and describe qualitatively the existence of such syntax correlations and reference locality in real programs. In Figure 2, the source and the assembly codes of a simple function *copy_disjunct* from

Parser of SPEC2000 are given. This function is invoked many times to build a new copy of a disjunct list. The second example *bsW* is extracted from *Bzip* of SPEC2000 (Figure 3). This function is also invoked multiple times to perform bit-stream I/Os.

Fig. 2. Example I: source and assembly codes of function *copy_disjunct* from Parser

The store/load syntax correlation and reference locality can be observed in several program constructs.

Register Save and Restore in Procedures and Functions: As shown in Figure 2, store/load dependences can be established perfectly with a matching pair of the base register (*\$sp*) and displacement for saving and restoring register contents when the function *copy_disjunct* is invoked. Although the invocations of *xalloc* and *copy_connectors* may change the value of the *\$sp*, the original value in the *copy_disjunct* is restored after returning from the function calls.

Access Records in Linked Data Structures: In the same example, the pointers (*d*, *d1*) are used to copy and construct a new node in the target linked structure. Different records (also pointers in this case) in each node of the old and the new linked structures are accessed using pointers *d*, *d1*. In the assembly code, the two pointers are loaded in registers *\$s0*, *\$s1* and are used as the base registers to access these records with small variations of the displacement value. The syntax correlations and reference locality among these accesses are clearly demonstrated in the assembly code.

Access Array Variables: Similar store/load correlations are also observed in accessing array data structures in several studied workload. For example, intensive array accesses are observed in several functions in *Gcc* of SPEC2000. Nearby references to different elements of the same array with the same base address provide syntax correlated stores and loads.

Access Global Variables: As shown in Figure 3, three global variables, *bsBuff*, *bsLive* and *bytesOut* are accessed when the function *bsW* is invoked. Due to the limited registers, these variables are loaded/stored multiple times based on the same global pointer *\$gp*. The access of global variables exhibits both the syntax correlation and the spatial locality.

Fig. 3. Example II: function *bsW* from *Bzip*, (a) source code; (b) assembly code; (c) partial assembly code from caller *SendMTFValues*

Access Local Variables: In the *bsW*, the callee-saved registers *\$s0* and *\$s1* are freed up for local usages to avoid saving parameters of *n* and *v* from registers *\$a0* and *\$a1* to the local stack and retrieving them later for computations. However, in functions that involve more complex computations and/or more temporary local variables, it is inevitable to increase the local stack accesses using the stack pointer *\$sp* and/or the frame pointer *\$s8* that also display strong syntax correlations and spatial locality.

Save/Restore Base Registers: There are evidences that the syntax correlation is still hold even if the base register has been updated between two memory accesses. This is due mainly to the fact that a base register may be freed up for other usages and

the original base address is restored before the next memory reference. In Figure 3, we also show a partial assembly code from a caller *SendMTFValues* of the *bsW*. In this caller, *\$s1* is used as a base register before calling the *bsW*. After returning from the *bsW*, *\$s1* continues to be used as a base register. Although *\$s1* has been updated in the *bsW*, the original base address is restored to keep the syntax correlation alive.

IV. ESTABLISHING A SYMBOLIC CACHE

An SC is a small data cache which is addressed by the encoding content of load/store instructions. The SC can be accessed once loads/stores are fetched out of the instruction cache. As a result, pipeline stages involving register file access, address generation/translation, and cache access can be bypassed. The impact of pipeline performance using an SC is very similar to that of using the VF in memory renaming techniques [25], where the speculative load data is fetched out of the VF indirectly through a store/load correlation table. In this paper, we focus on the accuracy of load data speculation using the SC. We omit discussions of integrating the SC into a pipeline microarchitecture.

It is essential to properly extract the *symbolic address* from the encoding bits of load/store instructions to capture the syntax correlations. A typical memory instruction consists of an opcode, a register source/destination, and a memory source/destination. Intuitively, we can use the memory source/destination to form a 32-bit symbolic address as illustrated in Figure 4. The least significant 16 bits are extracted from the displacement value, and the base register ID (5 bits) are inserted next to the displacement. Although simple, this approach suffers aliasing problems because multiple memory addresses can be mapped to the same symbolic address. In addition, this simple symbolic address formation creates other access and alignment problems.

Fig. 4. Extracting symbolic address from memory instructions

- *Aliasing of Symbolic Address:* With the simple address mapping in Figure 4, a 32-bit memory address is represented by a 21-bit symbolic address. Therefore, multiple memory addresses can be expressed by the same symbolic address. An obvious example can be found in stack accesses for local variables and for saving and restoring registers during procedure/function calls. Although accessing a different stack frame in each procedure invocation, the same stack pointer (*\$sp*) and frame pointer (*\$s8*) with a small range of displacement values are commonly used. The contents in the SC for local variables and saved registers are likely overwritten in the callee procedures and cannot be reused after returning from the procedures.
- *Uneven SC Index Distribution:* It is well-known that displacement values in memory references are unevenly distributed with a high percentage of ‘0’ and a few other constants. Using a portion of the high-order displacement bits as the index to the SC may potentially generate heavy conflict misses.
- *Word/Byte Alignment:* The most difficult problem lies in the difference of the line boundary between a symbolic

and a L_1 cache lines. This alignment problem is due to the fact that offset bits of a cache line are not always the same between the symbolic and the real addresses. It is essential to properly align the data layout in the symbolic cache according to the symbolic address to capture the spatial locality of memory references.

A. Procedure Coloring and Index Randomization

In order to alleviate the stack access aliasing problem in different procedures, various procedure coloring techniques can be constructed. A straight-forward technique is to maintain a global counter called *P-color*. The P-color is incremented whenever a procedure call is encountered. It is decremented after returning from a procedure. The P-color can be incremented contiguously in nested or recursive procedures before being decremented. Stack accesses between a caller and its callees can be differentiated by the P-color to avoid conflicts in the SC.

The P-color can be concatenated with the symbolic address for stack accesses. The width of the P-color counter is flexible. Figure 5 (a) illustrates the symbolic address after adding a 6-bit P-color. It is important to know that the P-color is only applied to stack accesses which use *\$sp* and *\$s8* as the base register. Other memory accesses do not add the P-color to allow sharing of global variables among different procedures or functions.

Fig. 5. (a) Adding procedure color to symbolic address; (b) Index randomization in accessing the SC

An uneven distribution of the index bits extracted directly from the displacement value has a potential to create heavy conflict misses in the SC. This problem comes from the fact that high-order displacement bits are often all zeros and can be dealt with by a simple randomization technique. Instead of extracting index bits from the symbolic address directly, randomized index bits can be formed by *exclusive-ORing* the original index bits from the displacement with the bits from the base register ID and the P-color as illustrated in Figure 5 (b). In this example, it is assumed that the SC has 64 sets with 64-byte line size. The six index bits are obtained by *exclusive-ORing* normal index bits in position 6 to 11 with the base register ID and partial P-color bits starting at position 16 through 21.

B. Word/Byte Alignment

One remaining issue is the data alignment between the SC and the L_1 data cache. The symbolic address within a cache line, i.e. the last few offset bits, may not be the same as the offset bits in the real address. In order to exploit spatial reference locality, the cache line fetched from L_1 needs to be rearranged in the SC such that the data layout can be aligned with the symbolic address. The basic alignment algorithm works as follows. When a memory request misses the SC, the target cache line is fetched from the memory hierarchy and loaded into the SC. The target byte/word is placed in the SC according to offset bits of the symbolic address. For example, assume there are eight access units in a cache line as shown in Figure 6. The symbolic offset of the target unit is *010* while the offset of the real address is *101*. In this case, the target data *101* is loaded into unit

010 in the SC. The remaining units are loaded according to the location of the target unit. There are thus two important aspects to consider for a proper data alignment:

Fig. 6. Data alignment in symbolic cache

- *Granularity of Data Alignment:* Depending on memory access granularity, it is conceivable that the data alignment can be performed at byte, half-word, word, or double-word level. The byte-level alignment can accommodate accesses by other granularity with the expense of maintaining more valid bits for the alignment information.
- *Handling Underflow/Overflow Data:* Since the line boundaries of the SC and the L_1 caches may be different, only a partial line can be filled on each SC miss. In addition, there is excessive data from the target L_1 cache line that cannot fit into the requested line location in the SC. The simplest and most natural solution is to only fill a partial SC line and drop the unfitted data. Other options include fetching two adjacent L_1 lines for each requested SC line, and/or to search and place the overflow L_1 data into the correct second SC line.

Performance evaluation on these design options will be given in the next section. It is important to keep the SC design simple since the primary goal of establishing the SC is to provide a zero-cycle load.

V. PERFORMANCE EVALUATION

Performance evaluations of three load value speculation methods are given including the last-value and stride-based value prediction (*VP*), the memory renaming (*MR*), and the proposed symbolic cache (*SC*). Our primary focus is to compare the prediction accuracy among these three mechanisms. All simulations are carried out on the *Sim-Save* model of SimpleScalar. Twelve integer programs, *Go*, *Li*, *M88k*, *Perl* from SPEC95 and *Bzip*, *Gcc*, *Gzip*, *Mcf*, *Parser*, *Twolf*, *Vortex*, *Vpr* from SPEC2000 are used. Version 2.7.2.3 *ssbig-na-sstrix-gcc* compiler with options: (*-funroll-loops -O2*) is used to generate the binary code. For each workload, we skip the first 900 million instructions, use the next 100 million instructions to warm up the caches and tables, then collect simulation statistics from the next 500 million instructions.

A. Data Alignment

We first investigate and evaluate different alignment granularity. Table I shows matches of the least-significant two bits between the symbolic and the real addresses with different memory access granularity in the simulated programs. On the average, 87.4%, 3.2% and 9.4% of memory references are accessing word, half-word, and byte respectively. Mismatches of the two bits for the three access granularities are about 0%, 0.5% and 4.5%. The word access is always aligned at the word boundary for both the real and the symbolic addresses. On the other hand, the word alignment creates 5% of mismatches for half-word and byte accesses. Since the word alignment reduces extra valid bits significantly, we will simulate both byte and word alignments and show their impact on the SC accuracy.

TABLE I

MATCHING OF THE TWO LEAST-SIGNIFICANT ADDRESS BITS BETWEEN REAL AND SYMBOLIC ADDRESSES FOR ACCESSING WORD, HALF-WORD AND BYTE

With regards to the line-fill on SC misses, preliminary studies show that the option of filling the entire SC line by fetching potentially more than one L_1 cache lines provides very limited benefit. Moreover, to place the entire target L_1 line into the SC on each miss does not benefit the accuracy much either. Therefore, only the simple partial SC line-fill by dropping any unfitted data is considered in subsequent evaluations.

B. Sensitivity of P-color and Index Randomization

Table II shows the accuracy of load value speculation using a 4KB SC with the word-alignment and 0, 2, 4 P-color bits. In general, we observe an average improvement from 68.8% to 70.4% by adding a 2-bit P-color. A few benchmark programs show no improvement at all with the simple P-color mechanism. After examining dynamic function calls in these programs, we found that there are very few nested calls and the program execution tends not to frequently traverse back and forth among multiple procedure levels. For instance, in Gzip, about 98% of the calls are labeled at level 6. We also observe that there is no benefit in increasing the number of bits in the P-color. With more P-colors, more levels of procedure invocation can be differentiated. However, analysis of application programs reveals that perfectly-nested or deeply-recursive procedures that benefit with more P-colors rarely exist. The actual execution path normally traverses among a few levels of procedures. Also, due to a small SC, the data from ancient ancestors is difficult to hold anyway.

TABLE II

LOAD ACCURACY USING THE SC WITH/WITHOUT THE P-COLOR

The benefit of index randomization is more evident in Table III, in which the accuracies of three 4KB SC configurations are displayed. By randomizing the index, a 4-way set-associative SC can achieve the accuracy approaching to that of a fully-associative SC. On the other hand, without this process, it degrades the accuracy of the 4-way design from 70.4% to 64.9%.

TABLE III

LOAD ACCURACY USING THE SC WITH INDEX RANDOMIZATION

These results suggest that the effective working set between base register updates is very small. Once the content of a base register changes, the old data in SC based on the same base register becomes stale. Because the original index bits are likely to be all zeros (Figure 5), stores and loads using the same base register may locate in very few sets even with index randomizations. Given the fact that the randomized 4-way SC achieves an accuracy comparable to that of a fully-associative SC, the

4 lines in each set are enough to hold the working set for each base register ID. Although higher set associativities increase the capacity in each set to hold more lines for each base register, frequent updates of base registers wipe out the corresponding correlated data in the SC.

C. Comparison of Three Data Speculation Methods

The accuracies of three load value speculation mechanisms are evaluated. Both byte and word alignments for placing a line in the SC are considered. Also, index randomizations and a 2-bit P-color are applied to improve the load accuracy. For a fair comparison, we simulate the three methods using comparable hardware with respect to the extra storage requirement to build additional tables and caches.

The VP scheme establishes a value history table to remember the recent value of each load. For matching the PC of a load, proper tags are maintained in the value history table. In addition, an increment value is needed in each entry to accommodate a stride-based predictor. The MR scheme uses a Value File (VF) to keep store/load correlated values for later accesses. In addition, two extra tables are needed. The Store/Load Cache (SLC) saves pointers to the VF. The SLC is addressed by the PCs of loads and stores with tags for matching the correct PC for indirect accesses to the VF. The Store-Address Cache (SAC) also records pointers to the VF. The SAC is accessed by load/store addresses for establishing load/store correlations. Again, address tags are necessary to make a correct correlation. The SC is simply a data cache addressed by the symbolic address. There is no extra hardware except for a small tag array in which each tag along with a few valid bits is associated with a 64-byte symbolic cache line.

We consider six configurations for accuracy comparisons as shown in Table IV. The hardware requirement is represented by the total number of entries in the respective tables and caches. Because of the additional tag arrays, the storage requirement for the VP and the MR are actually about 40-50% and 10-15% more than that of the SC in each configuration. Note that in this first-cut estimation, extra control logic is not considered.

TABLE IV

SIX CONFIGURATIONS FOR ACCURACY COMPARISONS

Figure 7 plots the average accuracy curves based on the twelve integer programs for the three data speculation methods. Generally speaking, the SC has the highest accuracy, especially with small configurations. For example, more than 70% of the loads can obtain correct values from a small 4KB SC. These results demonstrate the existence of store/load syntax correlations and spatial locality that can be captured effectively by small SCs. The MR scheme, on the other hand, requires 8 times of the hardware storage to reach about 67% accuracy. The MR scheme performs poorly with small configurations primarily because of misses to the small SLC/SAC for establishing correct store/load correlations. In addition, the correlation must be established before a correct value can be obtained. The MR scheme shows more improvement when the configuration size increases. With bigger SLC/SAC, data dependence links can be

Fig. 7. Average accuracies of three data speculation methods

built more precisely than those approximated by the symbolic address. However, the SC still maintains an edge by capturing the spatial locality. The last/stride value predictor generally has the worst accuracy. The accuracy improvement is leveling off with larger value history tables. This confirms a poor correlation between the load value and its instruction address.

The byte alignment does not improve the accuracy much. For a 4KB SC, for instance, the byte alignment improves the average accuracy of the word alignment from 70.4% to 71.1%. As shown in Table I, there is very little or no difference between byte or word alignment for a majority of the programs. The two programs that benefit the byte alignment the most are *Bzip* and *Gzip* because of their high percentage of sub-word accesses and mismatches of the least-significant 2 bits between real and symbolic addresses.

The SC size plays a minor role in providing accurate load values. Again, this is due to the fact that the working set between base register updates is very small. Since the randomized SC index is still mapped to very few sets for each base register, increasing the SC size (i.e. the number of sets) does not improve the capacity for loads using a specific base register.

Now considering the third configuration with a 4KB SC, the average prediction accuracies are 55.0%, 56.6%, 70.4% and 71.1% for the VP, the MR, and the SC with word alignment (SC-word) and the SC with byte alignment (SC-byte) respectively as shown in Figure 8. Among the twelve integer programs, *M88k*, *Perl*, and *Gcc* show very good syntax correlations with over 80% of prediction accuracies, while *Li*, *Gzip*, *Twolf*, *Vortex*, and *Vpr* show reasonable accuracies over 70%. *Go*, *Bzip*, *Mcf* and *Parser*, on the other hand, have poor accuracy, especially for *Go* with an accuracy only about 47%. Recall that in order to hold the syntax correlation, the base register content must remain unchanged between two correlative memory instructions. We found out in *Go*, about 64% of the loads are executed using a newly updated base register. On the other hand, only 22% and 24% respectively for the loads in *Gcc* and *M88k* are executed right after their base registers have updated. More detailed analysis with respect to the base register updates will be given in the next Section V-D.

The SC scheme does not perform well against the other two schemes under *Bzip* and *Parser*. In *Bzip*, a main function *full-GtU* that finds matches of character strings, has shown good value locality and good dynamic store/load correlations established by the MR scheme. However, the SC handles this function poorly because the base addresses of the matching strings are calculated right before loading characters from the two strings. A similar behavior has also found in *Parser*.

Fig. 8. Accuracy of three data speculation methods for individual programs (based on configuration 3)

In Figure 9, we break down correct and incorrect load value speculations using the SC with respect to the base register IDs. We separate base registers into 5 groups: *\$v*, *\$a*, *\$s+\$t*, *\$gp* and *\$sp+\$s8*, each represents 20.5%, 26.4%, 11.3%, 18.2% and 22.5% of the total loads, respectively. (Note there is about 1%

of the loads using other registers.) The accuracies of the 5 base register groups are 29%, 73%, 63%, 98%, and 94%. As expected, it is highly accurate to access global variables and local stack frames. For other loads, the compiler first picks *\$v* and *\$a* as temporary registers to hold base addresses for memory accesses. The base address is often computed or loaded from memory for an indirect access right before the load that results in an incorrect values from the SC. The *\$a* registers, which show higher accuracy, are also used for passing parameters to callee functions. We observe that many functions have memory addresses (pointers) as parameters that are passing through the *\$a* registers. In each callee function, the *\$a* registers are frequently used as a base without any modification. We also found in *Gcc* that certain memory addresses are passing through several function levels using the *\$a* registers. Thus, memory loads based on *\$a* can potentially keep the correlations alive through several function levels.

Fig. 9. Correct / incorrect load value speculations with respect to different base register groups

D. Accuracy Regarding Base Register Update

The syntax correlation holds when the content of the base register remains unchanged from the last memory reference with the same symbolic address. Figure 10 shows the average accuracy of all the loads with respect to the distance to the last update of the base register. For example, the distance is equal to 1 for a load when the base register of the load is used for the first time as a base register after an update to the register. Similarly, the distance is equal to 2 if a register is used for the second time as a base register for either load or store after the content of the register has updated. The distances of 20 or longer are represented by a single data point. In general, the accuracy goes up with the distance due to the locality of references. A cold miss is encountered when the distance is equal to 1 unless the latest update did not change the content of the base register from the previous use of the same base register.

Fig. 10. Load accuracy and distribution with respect to the distance to the last base register update

A few observations can be made from the figure. First, when the distance is 3 or longer, the speculative load data from the SC is very accurate with an average accuracy about 98%. This indicates a very strong reference locality based on the symbolic addresses of nearby stores and loads.

Second, instead of all cold misses, the average accuracy is 36% when the distance is equal to 1. This accuracy comes from restoring base register content before the load. Unfortunately, a significant portion (39%) of the loads use a base register at the first time after its updates. With only 36% of accuracy, these loads produce 25% inaccurate data with respect to the total loads. Therefore, the distance-1 loads are the major factor for the overall accuracy. For example, in the two high-accuracy programs, *Gcc* and *M88k*, only 22% and 24% of loads are distance-1 with an accuracy of 46% and 68%, respectively. On the other hand, *Go* has 64% of loads are distance-1 with a poor accuracy of 18%.

Third, about 24% of the loads have distances of 20 or longer. This long distance comes mainly from access global variables, also from some local variable accesses. An average accuracy of 98.4% is obtained for these long-distance loads.

Compiler optimization techniques may be applied to improve the syntax correlations of stores/loads. For example, we observe that parameters are sometimes passed to the callee through the caller's stack frame. Accessing the parameters before an update of the frame pointer may keep the correlation alive. Further discussions in this direction is out of the scope of this paper.

VI. CONCLUSION

A new load data speculation method, based on instruction syntax correlations of stores and loads, has been introduced in this paper. Instead of establishing the store/load correlation dynamically at runtime, the proposed method establishes a small symbolic cache to capture existing syntax correlations and memory reference locality. The symbolic cache is addressed by the encoding content of store/load instructions to enable data accesses in the front-end of the processor pipeline to shorten load-to-use latency. Performance evaluation of SPEC integer programs has demonstrated that the proposed method can achieve an accuracy over 70% with a small 4KB symbolic cache. With compiler helps to reduce base register updates and to better utilize displacement values, further improvement of the SC accuracy may still be possible.

Acknowledgment:

This work is supported in part by NSF grants MIP-9624498, EIA-0073473 and by Intel research and equipment donations. Anonymous referees provide very helpful comments.

REFERENCES

- [1] T. Austin and G. Sohi, "Zero-cycle loads: microarchitecture support for reducing load latency", *Proc. of 28th Int'l Symp. on Microarchitecture*, Ann Arbor, MI, Dec. 1995, pp. 82–92.
- [2] M. Bekerman, S. Jourdan, R. Ronen, G. Kirshenboim, L. Rappoport, A. Yoaz, and U. Weiser, "Correlated Load-Address Predictors," *Proc. of 26th Int'l Symp. on Computer Architecture*, Atlanta, GA, May 1999, pp. 54–63.
- [3] M. Bekerman, A. Yoaz, F. Gabbay, S. Jourdan, M. Kalaev, and R. Ronen, "Early Load Address Resolution Via Register Tracking," *Proc. of 27th Int'l Symp. on Computer Architecture*, Vancouver, Canada, June 2000, pp. 306–315.
- [4] D. Burger and T. Austin, "The SimpleScalar Tool Set, Version 2.0", Technical Report #1342, CS Department, University of Wisconsin-Madison, June 1997.
- [5] B. Calder, G. Reinman, and D. Tullsen, "Selective Value Prediction," *Proc. of 26th Int'l Symp. on Computer Architecture*, Atlanta, GA, May 1999, pp. 64–75.
- [6] C. Chen and A. Wu, "Microarchitecture Support for Improving the Performance of Load Target Prediction," *Proc. of 30th Int'l Symp. on Microarchitecture*, Triangle Park, NC, Dec. 1997, pp. 228–234.
- [7] B. Cheng, D. Connors, and W. Hwu, "Compiler-Directed Early Load-Address Generation," *Proc. of 31st Int'l Symp. on Microarchitecture*, Dallas, TX, Dec. 1998, pp. 138–147.
- [8] B. Chung, J. Zhang, J-K. Peir, S. Lai, and K. Lai, "Direct Load: Dependence-Linked Dataflow Resolution of Load Address and Cache Coordinate," *Proc. of 34th Int'l Symp. on Microarchitecture*, Austin, TX, Dec. 2001, pp. 76–87.
- [9] R. Eickemeyer and S. Vassiliadis, "A Load-Instruction Unit For Pipelined Processors," *IBM Journal of Research and Development*, Vol. 37(4), pp. 547–564, July 1993.

- [10] J. Gonzalez, and A. Gonzalez, "Speculative Execution via Address Prediction and Data Prefetching," *Proc. of 1997 Int'l Conf. on Supercomputing*, Vienna, Austria, Aug. 1997, pp. 196–203.
- [11] J. Gonzalez, and A. Gonzalez, "The Potential of Data Value Speculation to Boost ILP," *Proc. of 1998 Int'l Conf. on Supercomputing*, Melbourne, Australia, June, 1998, pp. 21–28.
- [12] T. Horel and G. Lauterbach, "UltraSPARC-III: Designing Third-Generation 64-Bit Performance", *IEEE Micro*, May/June 1999, pp. 73–85.
- [13] K. Hua, A. Hunt, L. Liu, J-K. Peir, D. Pruett, and J. Temple, "Early Resolution of Address Translation in Cache Design," *Proc. of 1990 Int'l Conf. on Computer Design*, Boston, MA, Sep. 1990, pp. 408–412.
- [14] R. Kessler, "The Alpha 21264 Microprocessor," *IEEE Micro*, Vol. 19(2), March/April 1999, pp. 24–36.
- [15] M. Lipasti, C. Wilkerson and J. Shen, "Value Locality and Load Value Prediction", *Proc. of the 7th Int'l Conf. on Architectural Support for Programming Languages and Operating Systems*, Boston, MA, Oct. 1996, pp. 138–147.
- [16] M. Lipasti and J. Shen, "Exceeding the Limit via Value Prediction," *Proc. of 29th Int'l Symp. on Microarchitecture*, Paris, France, Dec. 1996, pp. 226–237.
- [17] W. Lynch, G. Lauterbach and J. Chamdani, "Low Load Latency through Sum-Addressed Memory (SAM)," *Proc. of 25th Int'l Symp. on Computer Architecture*, Barcelona, Spain, June 1998, pp. 369–379.
- [18] Q. Ma, J-K. Peir, L. Peng, and K. Lai, "Symbolic Cache: Fast Memory Access Based on Program Syntax Correlation of Loads and Stores," *Proc. of 2001 Int'l Conf. on Computer Design*, Austin, TX, Sep. 2001, pp. 54–61.
- [19] A. Moshovos and G. Sohi, "Streamlining Inter-Operation Memory Communication via Data Dependence Prediction," *Proc. of 30th Int'l Symp. on Microarchitecture*, Triangle Park, NC, Dec. 1997, pp. 235–245.
- [20] A. Moshovos and G. Sohi, "Read-After-Read Memory Dependence Prediction," *Proc. of 32nd Int'l Symp. on Microarchitecture*, Haifa, Israel, Nov. 1999, pp. 177–185.
- [21] D. Papworth, "Tuning the Pentium Pro Microarchitecture," *IEEE Micro*, Vol. 16(2), April 1996, pp. 8–15.
- [22] Y. Sazeides and J. Smith, "The Predictability of Data Values," *Proc. of 30th Int'l Symp. on Microarchitecture*, Triangle Park, NC, Dec. 1997, pp. 248–258.
- [23] T. Slegel, et al., "IBM's S/390 G5 Microprocessor Design," *IEEE Micro*, Vol. 19(2), March/April 1999, pp. 12–23.
- [24] P. Song, "IBM's Power3 to Replace P2SC," *Microprocessor Report*, Vol. 11(15), Nov. 1997, pp. 1–11.
- [25] G. Tyson and T. Austin, "Improving the Accuracy and Performance of Memory Communication Through Renaming," *Proc. of 30th Int'l Symp. on Microarchitecture*, Triangle Park, NC, Dec. 1997, pp. 218–227.
- [26] K. Wang, and M. Franklin, "Highly Accurate Data Value Prediction using Hybrid Predictors," *Proc. of 30th Int'l Symp. on Microarchitecture*, Triangle Park, NC, Dec. 1997, pp. 281–290.

PLACE
PHOTO
HERE

Lu Peng received his B.E. and M.S. degrees in Computer Science from Shanghai Jiaotong University, China. He is currently a PhD student with the Computer and Information Science and Engineering department at University of Florida. His research interests include computer architecture, distributed system and computer networks. He has been a student member of the ACM since 1999.

PLACE
PHOTO
HERE

Jih-Kwon Peir received his B.S. degree in Engineering from Cheng-Kung University, Taiwan, M.S. degree from University of Wisconsin-Milwaukee, and Ph.D. degree from University of Illinois, Urbana-Champaign, both in computer science. He worked at IBM T.J. Watson Research Center (1986-1992), as a research staff member involved in mainframe processor designs. During 1992-93, he was the deputy director of Computer Technology in the Industrial Technology Research Institute in Taiwan in charge of an Intel Pentium-based SMP development project.

Since 1995, he spent several summers visiting Intel's Microprocessor Research Lab and IBM's Almaden Research Center. He is currently an Associate Professor in the Computer and Information Science and Engineering Department at University of Florida. Dr. Peir's research interests include computer system

architectures, designs, and performance evaluation. He received an IBM Faculty Development Partnership Award in 1995 and an NSF Faculty Early Career Development Award in 1996. He was a co-author of two best paper awards (1990, 2001) in IEEE-ICCD conference. He serves as a subject area editor of the *Journal of Parallel and Distributed Computing* and is on the editorial board of the *IEEE Transactions of Parallel and Distributed Systems*.

PLACE
PHOTO
HERE

Qianrong Ma received his B.E. and M.E. degrees in electrical engineering from Tsinghua University, China, in 1993 and 1996, respectively. He spent one year with Stone Corporation, China, as a product development engineer. From August 1997 to May 2000, he studied in the Computer and Information Science and Engineering Department, University of Florida in the area of computer architecture. He earned his M.S. degree in Computer Engineering in May 2000. He is currently a senior member of technical staff at Server Technology Division, Oracle Corporation.

PLACE
PHOTO
HERE

Konrad Lai received the B.S. degree in electrical engineering from Princeton University, Princeton, NJ, in 1976, and the M.S. degree in computer engineering from Carnegie Mellon University, Pittsburgh, PA, in 1978. He is currently a Principal Researcher and a manager in Microprocessor Research, Intel Labs, in Hillsboro, OR. He has been with Intel for over 20 years, working on microprocessor, memory, and system architecture. He holds 25 patents and has authored or co-authored eight papers. Mr. Konrad is a Member of ACM and the IEEE Computer Society.