

Techniques for Fast Instruction Cache Performance Evaluation

DAVID B. WHALLEY

Department of Computer Science B-173, Florida State University, Tallahassee, FL 32306, U.S.A.

SUMMARY

Cache performance has become a very crucial factor in the overall system performance of machines. Effective analysis of a cache design requires the evaluation of the performance of the cache for typical programs that are to be executed on the machine. Recent attempts to reduce the time required for such evaluations either result in a loss of accuracy or require an initial pass by a filter to reduce the length of the trace. This paper evaluates techniques that attempt to overcome these problems for instruction cache performance evaluation. For each technique variations with and without periodic context switches are examined. Information calculated during the compilation is used to reduce the number of references in the trace. Thus, in effect references are stripped before the initial trace is generated. These techniques are shown to significantly reduce the time required for evaluating instruction caches with no loss of accuracy.

KEY WORDS: Instruction Cache Cache Simulation Trace Generation Trace Analysis

INTRODUCTION

The time required to generate and analyze trace data is proportional to the number of references in the trace¹. Since realistic program traces can be quite lengthy, a trace is often only collected from a portion of the program's execution. However, it has been shown that the cache performance can vary greatly in different portions of a program's execution². Cache performance measurements obtained when unrealistic input data is used to shorten the length of the trace would also be of questionable value.

Most recently designed machines store instructions in a separate cache from data. There are many advantages to having an instruction-only cache. Most machines do not allow modification of instructions during the execution of a program. By only allowing read accesses to a cache, the design of an instruction cache is simplified. Also, instruction references typically have higher locality than data references. By using separate caches, the designers of a system can offer different configurations for each cache which may improve the overall performance. Since the type of reference, instruction or data, issued to the memory system is known by the CPU, there can be separate ports for instructions and data. Thus, the bandwidth between cache memory and the CPU can be improved.

Since instruction caches have become more common, there has been much recent work that attempts to reorganize the code of programs to improve the instruction cache performance³⁻⁵. In fact, most compiler

optimizations can affect instruction cache performance since the optimizations can change the order and number of instructions that are executed by a program. Since these optimizations can affect different portions of a program's execution, the effect of a compiler optimization on instruction cache performance can only be accurately evaluated if the complete trace of instruction addresses from the program's execution is used.

This paper describes techniques that use compile-time analysis to reduce the time required for evaluating instruction cache performance. Unlike many data references, the address of each instruction remains the same during the entire execution of a program. Information can therefore be calculated prior to the instruction cache simulation that can be used to reduce the number of references that the cache simulator needs to process. Thus, in effect address references are being stripped before the initial trace is generated.

RELATED WORK

Traditionally, the problem of evaluating the cache performance during the execution of a program has been separated into two tasks, trace generation and trace analysis. The first task is to generate the trace of addresses that will be presented to the cache. The address trace is typically either written to secondary storage, such as a file on disk, or stored in a trace buffer in memory. The second task is to analyze the addresses that were generated. This analysis is usually accomplished through the use of a separate program that reads the generated address trace and simulates the behavior of the cache.

Trace-driven cache simulation has long been the primary method used to analyze the performance of a cache. Simple approaches for generating trace data and simulating caches, however, can be both very time and space consuming. Two common methods used for generating trace data are forcing a program to trap after the execution of each instruction or to record references while simulating the execution of each instruction. Each of these methods can result in a program executing a 1000 times slower than normal execution⁶⁻⁸. Trace data is typically stored in secondary storage and then later read by a cache simulator that will perform the analysis of the data. Realistic trace data, however, requires at least several million references which may make infeasible the use of disk as the storage media². Therefore, there has been much work on the problem of reducing the

space and time requirements for trace-driven simulation.

There are other trace generation methods that are much faster than simulation or trapping. One method is to modify the microcode of a microprogrammed machine. Relative to techniques using traps or simulated program execution, this method does not impose a large run-time penalty⁹. The microcode of a machine, however, is often not accessible to the typical user. Even when it is accessible, it often requires great expertise to modify without adversely affecting the operation of the machine. Also, modification of microcode would not be applicable for machines which are not microprogrammed (e.g. many RISCs). Another method is to use a hardware monitor¹⁰ which can collect address references without perturbation of a machine. These monitors are specialized and expensive hardware devices and are also not available to a typical user. In addition, a monitor that watches a bus to obtain the addresses of references would not be appropriate with an on-chip instruction cache.

Program instrumentation, or inline tracing, is a technique that requires little overhead for generating a trace of addresses^{2,8,11}. Instructions are inserted to record addresses during the program's execution and must not change the normal execution behavior of the program. Therefore, the values in data registers or the condition codes may have to be saved and restored. The trace generation overhead can be reduced by generating a subset of the trace from which an entire trace can be regenerated. One can also optimize the placement of the instrumentation code to produce the reduced trace with respect to a weighting of the control-flow graph¹².

Unfortunately, even with fast trace generation techniques, evaluating the cache performance of a program's execution can be quite time-consuming since the largest factor in the time required for cache performance evaluation is analyzing the trace of addresses. Even after tuning a cache simulator for a specific cache configuration, a cache simulator can still require an order of magnitude more time than generating the trace itself². Therefore, there has been much attention given to reducing the number of references that need to be traced.

There have been several methods proposed to improve cache simulation times by reducing the number of references in the trace. One widely known method for reducing the length of the trace data is called trace stripping¹³. This approach first simulates a direct-mapped cache and records only the references that are

misses since hits do not result in changes to the cache state. The reduced trace can then be used to simulate caches with a greater number of sets and associativity as long as the line size is not changed and context switches are not introduced. Unfortunately, this technique requires the entire trace to be processed by the cache simulator once. Furthermore, the reduced trace may still be quite lengthy, which can result in large files and slow simulations. There have also been several methods that allow different cache configurations to be evaluated during a single simulation^{14,15}.

Another method to avoid processing the entire trace of references is to instead use several discrete samples of traces from the program execution to predict cache performance measures¹⁶. While a significant loss of accuracy may not occur, the method may not have the desirable accuracy for measuring the effect of a new compiler optimization or reorganization technique. For instance, the optimization may be applied on a section of code that when executed is not in the samples of references that are collected.

On-the-fly analysis is a technique that avoids the I/O associated with storing the generated trace and retrieving the trace for input to the cache simulator. In this approach either the cache simulator is a separate process that reads a trace buffer containing the trace² or the cache simulator is linked directly to the program and trace information is received as arguments via function calls¹⁷. Even though the space and I/O requirements are diminished, the trace analysis can still be quite time-consuming since the entire trace is being processed.

TECHNIQUES FOR REDUCING INSTRUCTION CACHE EVALUATION TIMES

An optimizer of a compiler system was modified to be able to generate and analyze trace data. The program being measured is compiled in two passes. The first compilation pass serves to determine the address of each instruction. Rather than using a table lookup method, which may be rather complicated for highly encoded CISC architectures, a label is inserted before and after each instruction. A call is inserted at the beginning of the main function to invoke a routine that will dump out the size of each instruction, the difference between each

pair of labels.* Even the size of branch instructions, which may vary depending upon the distance to the branch target, may be accurately determined since no other trace instructions are inserted at this point. The first compilation pass also stores information about each function to be used by some of the techniques described later in the paper.

The second compilation pass inserts the instructions to invoke the cache simulator during the program's execution. The values of scratch registers and condition codes may have to be saved prior to the inserted instructions and restored afterwards. By using the data-flow information already calculated and available in the optimizer, these saves and restores are only inserted when necessary. A call to print the cache performance report is also inserted before any return instructions in the main function or calls to the `exit` function anywhere in the program. The cache simulator is linked with the program to allow on-the-fly trace analysis to occur while the program is being executed.

The techniques used to reduce the number of cache references processed by the cache simulator are described in the following sections. The cache references not processed are those references which are ascertained to be cache hits and will not change the state of the cache. Examples associated with the techniques are given in Motorola 68020 assembly code. Each technique builds upon the previously presented techniques. Technique A is a straight-forward approach. Techniques B and C recognize spatial locality to reduce the length of the trace and methods similar to these have been used in previous studies^{11,18}. Techniques D-G use cache configuration information and the control flow of the program to avoid processing additional references due to both spatial and temporal locality. Two variations of each of the techniques are discussed. The first variation allows periodic context switches and the second variation assumes there are no context switches.[†]

Technique A

For Technique A, a call to a trace routine is inserted before each basic block in the program. Note, that basic blocks for tracing are delimited by labels and branching instructions, including calls. A block number,

* On some older machines, such as the CDC 6600, the insertion of a label can change the address of an instruction since all branch targets had to be aligned on a word boundary. This author is not aware of any current machines with this characteristic. Furthermore, the current trend is to have instructions that are all the same size to simplify decoding.

that uniquely identifies the basic block being executed, is passed to the trace routine which uses the number to access information associated with the block.[‡] The trace routine in turn interfaces with the cache simulator, passing it the address and size of each instruction within the block. If the cache simulator was compiled to allow the option of periodic context switches, then a simulated context switch invalidating the entire instruction cache can occur between two instructions in the same basic block. Assembly code with inserted instructions for Technique A is given in Figure 1. The source code for the trace routine is given in Figure 2.

```

...
L142:  pea    #145          /* push block number 145 */
      jbsr  _traceblknum /* call trace routine */
      addq1 #4,a7         /* adjust stack pointer */
...
      jne   L67
      pea  #146          /* push block number 146 */
      jbsr  _traceblknum /* call trace routine */
      addq1 #4,a7         /* adjust stack pointer */
...

```

Figure 1: Assembly Code with Technique A

```

void traceblknum(blk)
int blk;
{
    register int i;

    /* Invoke the cache simulator for each instruction in the block. */
    for (i = blkinfo[blk].first;
         i < blkinfo[blk].first+blkinfo[blk].numinst; i++)
        cachesim(instsize[i], instaddr[i]);
}

```

Figure 2: Trace Routine for Technique A

[†] An earlier study by this author only evaluated techniques that allow periodic context switches¹⁹.

[‡] Though a *cache line* is sometimes referred to as a *block*, in this paper the term *block* indicates a *basic block* in the compiled program.

Technique B

The instructions inserted into the program for Technique B are identical to the instructions inserted for Technique A. The trace routine being invoked, however, interfaces with the cache simulator differently. If no context switches can occur during the execution of the block or the cache simulator was compiled to not allow context switches, then the cache simulator is invoked only once with the address of the first instruction within the basic block and the size of entire block passed as arguments. Thus, the cache simulator is treating the entire block as one large instruction. The actual number of references to the cache associated with the block, which is determined statically, may be greater than the number of instructions within the block. This situation occurs when one or more instructions span more than one cache line. The number of references, hits and misses, caused by invoking the simulator for the entire block as one large instruction is equal to the number of cache lines being referenced. The remaining references associated with the block due to spatial locality will always be cache hits and need not be simulated. Thus, the counter for the number of cache hits is incremented by the number of remaining references after the call to the simulator. If it is determined that a periodic context switch may occur during the execution of the block, then the cache simulator will be invoked for each individual instruction within the block as in Technique A. The source code for the trace routine with Technique B is given in Figure 3. The variables `c_hits` and `c_misses`, which represent the number of cache hits and misses, are global variables with initial values of zero. The portions that are in boldface are used only when periodic context switches are allowed.

Technique C

Technique C is similar to Technique B except for one difference. Instead of invoking the simulator once for an entire basic block, the technique invokes the simulator once for each sequence of executed blocks that are physically contiguous. The basic block number of the beginning of a sequence of blocks, also described as a superblock¹¹, is saved. A call instruction to the trace routine is inserted before any unconditional jump, call, or return instructions. Handling conditional transfers of control is a little more complicated. The trace routine should only be invoked when the conditional branch is taken. However, the conditional branch target block

```

void traceblknum(blk)
int blk;
{
    register struct blkinfo *p;
    int i;
    int j;

    /* If a context switch cannot occur during the execution of the
       block then invoke the cache simulator once for the entire block. */
    p = &blkinfo[blk];
    if (!(c_switch + p->worst > SWITCHTIME)) {
        i = c_hits + c_misses;
        cachesim(p->size, instaddr[p->first]);
        c_hits += (j = (p->refs - ((c_hits + c_misses) - i)));
        c_switch += j*HITTIME;
    }

    /* Else invoke the cache simulator for each instruction in the block. */
    else {
        for (i = p->first; i < p->first+p->numinst; i++)
            cachesim(instsize[i], instaddr[i]);
    }
}

```

Figure 3: Trace Routine for Technique B

could be in the middle of a different sequence of contiguous blocks. Therefore, the target of the conditional jump is replaced with the target of a newly created label. Assembly code is added at the end of the function that contains the new label, the call to the trace routine, the reset of the beginning block of a new sequence, and an unconditional jump back to the conditional jump's original target. An example of assembly code before and after modifications are made for Technique C is shown in Figure 4. The source code for the trace routine with Technique C is given in Figure 5. Again the portions that are in boldface are used only when periodic context switches are allowed.

Technique D

The remaining techniques use cache configuration information to reduce the number of references to be processed by the cache simulator. The compiler reads in information that indicates the line size and number of sets of a direct-mapped cache. The compiled program can be used to simulate caches with different

before	after
-----	-----
jbsr _foo	jbsr _foo
...	movl #15,_startblk /* follows a call inst */
jne L74	...
...	jne LN10 /* L74 was replaced */
jeq L78	...
...	jeq LN11 /* L78 was replaced */
jra L67	...
	pea #17 /* push last block in seq */
	jbsr _traceblknum /* call trace routine */
	addq #4,a7 /* adjust stack pointer */
	movl #32,_startblk /* reset start block */
	jra L67
	.
	.
	.
	LN10: pea #15 /* push last block in seq */
	jbsr _traceblknum /* call trace routine */
	addq #4,a7 /* adjust stack pointer */
	movl #22,_startblk /* reset start block */
	jra L74 /* jump to orig label */
	LN11: pea #16 /* push last block in seq */
	jbsr _traceblknum /* call trace routine */
	addq #4,a7 /* adjust stack pointer */
	movl #26,_startblk /* reset start block */
	jra L78 /* jump to orig label */

Figure 4: Assembly Code with Technique C

characteristics including greater set associativity. As in Puzak's trace stripping method¹³, if the number of sets is not decreased and the line size remains the same, then the program need not be recompiled.

The variation of Technique D that allows periodic context switches is accomplished as follows. An array, where an element is indexed by a unique block number, is used to indicate if each basic block in the program is currently within the cache. The compiler determines if each loop in a function, proceeding from the outermost loop first, can fit in the cache and does not contain any calls to other functions that are being measured. The same cache simulator used to analyze the cache performance during the execution of a measured program was linked with the compiler so it could be invoked to check if a loop will fit in the cache. If the loop does fit, then instructions are inserted in the preheader block of the loop to clear the array elements associated with the blocks

```

void traceblknum(blk)
int blk;
{
    register int i, refs, size;
    int j;
    register int worst;

    /* Sum information about the set of consecutive blocks. */
    refs = size = 0;
    worst = 0;
    for (i = startblk; i <= blk; i++) {
        refs += blkinfo[i].refs;
        size += blkinfo[i].size;
        worst += blkinfo[i].worst;
    }

    /* If a context switch cannot occur during the execution of the set
       of blocks then invoke the cache simulator once for the entire set. */
    if (!(c_switch + worst > SWITCHTIME)) {
        i = c_hits + c_misses;
        cachesim(size, instaddr[blkinfo[startblk].first]);
        c_hits += (j = (refs - ((c_hits + c_misses) - i)));
        c_switch += j*HITTIME;
    }

    /* Else invoke the cache simulator for each instruction in the set of blocks. */
    else {
        for (i = startblk; i <= blk; i++)
            for (j = blkinfo[i].first;
                 j < blkinfo[i].first+blkinfo[i].numinst; j++) {
                cachesim(instsize[j], instaddr[j]);
            }
    }
}

```

Figure 5: Trace Routine for Technique C

in the loop. Note, that if a preheader block is not available, then one will be created. For each block within the loop, instructions are inserted to determine if the block currently resides in the cache by checking the array element associated with the block. If the array element is currently set and a simulated context switch cannot occur while the instructions in the block are executed, then the number of cache hits and the context switch

information are adjusted.* It is assumed that the inserted instructions to perform these two checks are much less expensive than processing the reference by the cache simulator. If a context switch does occur, then the array elements within the loop are cleared. Thus, when a loop that fits in the cache is entered, except when a context switch occurs, the cache simulator is invoked at most once for each block in the loop. Technique C is used for code outside of loops or in loops that do not fit in the cache. For the example in Figure 6, which contains a loop with a single basic block, simulated context switches occur every 10,000 units of work. The context interval can be changed, but would require recompiling the program being measured.

```

        clr1    _blkmarker+952    /* clear blkmarker[238] */
        movl   #238,_lowmarker    /* first block cleared on context switch */
        movl   #238,_highmarker   /* last block cleared on context switch */
L405:   cmpl   #0,_blkmarker+952  /* check if block 238 is in cache */
        jeq    LN191              /* if not then have to invoke simulator */
        cmpl   #9992,_c_switch    /* check if context switch in 8 cycles */
        jge    LN191              /* if pending then invoke simulator */
        addq1  #8,_c_switch        /* adjust context switch information */
        addq1  #8,_c_hits         /* adjust number of cache hits */
LN192:  ...
        jne    L405
        movl   #-1,_highmarker    /* nothing to clear since not in loop */
        .
        .
        .
LN191:  movl   #1,_blkmarker+952   /* denote that block 238 is in cache */
        movl   #238,_startblk     /* first block in sequence */
        pea   #238                /* last block in sequence */
        jbsr  _traceblknum        /* call trace routine */
        addq1 #4,a7                /* adjust stack pointer */
        jra   LN192              /* will now execute insts in basic block */

```

Figure 6: Assembly Code with Technique D and Context Switches

The variation of Technique D that assumes that there are no context switches is similar to the one with context switches with the following differences. Each element of the array is used to represent the number of outstanding cache hits associated with a basic block. At the beginning of the execution of the program, each element is initialized to the negated number of cache references that would be required if the block was

* Note that adjusting this information does result in a stronger coupling between the trace generation and trace analysis tasks.

executed. An instruction is inserted at the beginning of each block to add the number of references to the associated array element. At the exit(s) of the loop, instructions are inserted to update the number of cache hits and misses for each block within the loop. First, the array element associated with a block in the loop is checked to determine if it was executed. If so, then the first execution of the block is simulated by the cache simulator, the remaining references are added to the number of cache hits, and the array element is reset to the negated number of cache references. Figure 7 shows the same example as in Figure 6, but with the assumption that no context switches can occur. Note that only one instruction is inserted with each block in the loop rather than six.

```

L405:  addq1    #8,_blkmarker+952    /* adjust number of hits for block 238 */
      ...
      jne     L405
      cmpl   #0,_blkmarker+952    /* check if block 238 was executed */
      jlt   LN191                 /* if not then skip over measurement code */
      movl   #238,_startblk       /* first block in sequence */
      pea   #238                  /* last block in sequence */
      jbsr   _traceblknum        /* call trace routine */
      addq1  #4,a7                 /* adjust stack pointer */
      movl   _blkmarker+952,d0    /* load remaining hits for block 238 */
      addl   d0,_c_hits           /* adjust number of cache hits */
      movl   #-8,_blkmarker+952  /* reset hits for block 238 */
LN191: ...

```

Figure 7: Assembly Code with Technique D and No Context Switches

Figure 8 shows the source code for the trace routine with Technique D. This trace routine is also used for Techniques E-G. Again the portions of the source code in boldface are used only when periodic context switches are allowed. The only difference between this trace routine and the one used for Technique C is that a set of markers are cleared when a periodic context switch occurred.

The cache configuration information read by the compiler is used to determine if the cache lines that are associated with a basic block are currently in the cache. If these cache lines are resident and each cache line was the last line to be referenced within its set, then there is typically no information that need be updated about the state of the cache in the simulator. The replacement algorithm for determining which cache line to replace within the set can be usage or non-usage based²⁰. Non-usage based algorithms, such as first-in-first-out

```

void traceblknum(blk)
int blk;
{
    register int i;
    register int refs, size;
    int j, k;
    register int worst;

    /* Sum information about the set of consecutive blocks. */
    refs = size = 0;
    worst = 0;
    for (i = startblk; i <= blk; i++) {
        refs += blkinfo[i].refs;
        size += blkinfo[i].size;
        worst += blkinfo[i].worst;
    }

    /* If a context switch cannot occur during the execution of the set
       of blocks then invoke the cache simulator once for the entire set. */
    if (!(c_switch + worst > SWITCHTIME)) {
        i = c_hits + c_misses;
        cachesim(size, instaddr[blkinfo[startblk].first]);
        c_hits += (j = (refs - ((c_hits + c_misses) - i)));
        c_switch += j*HITTIME;
    }

    /* Else invoke the cache simulator for each instruction in the set of blocks. */
    else {
        for (i = startblk; i <= blk; i++)
            for (j = blkinfo[i].first;
                 j < blkinfo[i].first+blkinfo[i].numinst; j++) {
                cachesim(instsize[j], instaddr[j]);

                /* Clear set of block markers if context switch occurred. */
                if (old_context_switch < context_switch) {
                    old_context_switch = context_switch;
                    for (k = lowmarker; k <= highmarker; k++)
                        blkmarker[k] = 0;
                }
            }
    }
}

```

Figure 8: Trace Routine for Technique D

(FIFO) or random, are not affected by the reuse of a line. The most common usage-based replacement algorithm for set-associative caches is least-recently-used (LRU)²¹. If the current cache line being referenced is also the last line that was referenced within the set, then LRU information need not be updated. Therefore, a program containing a loop that fits in a direct-mapped cache and processed using Technique D would not require recompilation when the number of sets or associativity is increased.

Technique D, as described so far, attempts to avoid calls to the cache simulator for blocks that will be in the cache already due to temporal locality. There are also situations when blocks will already be in the cache due to spatial locality. When a basic block in a loop that fits in the cache is totally contained in the last cache line referenced by the predecessor block and/or the first cache line of the successor block, then instructions are inserted to check the block markers of the predecessor and/or successor blocks. If set, then the call to the cache simulator is avoided and the context switch information and cache hit counter are adjusted.

Technique E

Technique E is similar to Technique D except that inter-procedural cache analysis is performed. Initially, a call graph of the functions being measured is constructed using the information provided from the first compilation pass. Then, it is determined if each subtree of the graph, a function and the routines that can be invoked from that function, can fit in the cache at the same time. If an entire function and the routines that can be invoked will fit in the cache, then Technique D, except for clearing the block markers, is used throughout the function. In a function that does not fit in cache, instructions are inserted preceding calls to a function being measured that does fit in the cache to clear the markers associated with that function. Loops with calls to functions that are being measured may also be candidates for avoiding references to be processed by the cache simulator. If the entire loop and the functions that can be invoked from the loop can fit in cache at the same time, then the block markers associated with the blocks in the loop and with each of the functions that can be invoked are cleared in the preheader block of the loop. Inter-procedural cache analysis can still be used when call graphs are cyclic (i.e. recursive). Technique E (and Techniques F and G), however, cannot be used with indirect calls since the function being invoked is not known at compile-time and the call graph cannot be accurately

constructed. The example given in Figure 6 for Technique D showed a loop with a single block that contained no calls. Figure 9 shows the same loop except with one call to a routine being measured. The two loop blocks and the three basic blocks in the routine fit in the cache at the same time.

```

        clr1    _blkmarker+952    /* clear blkmarker[238] */
        clr1    _blkmarker+956    /* clear blkmarker[239] */
        clr1    _blkmarker+1056   /* clear blkmarker[264] */
        clr1    _blkmarker+1060   /* clear blkmarker[265] */
        clr1    _blkmarker+1064   /* clear blkmarker[266] */
        movl    #238,_lowmarker   /* first block cleared on context switch */
        movl    #266,_highmarker  /* last block cleared on context switch */
L405:   cml    #0,_blkmarker+952  /* check if block 238 is in cache */
        jeq     LN191             /* if not then have to invoke simulator */
        cml    #9995,_c_switch    /* check if context switch in 5 cycles */
        jge     LN191             /* if pending then invoke simulator */
        addql   #5,_c_switch      /* adjust context switch information */
        addql   #5,_c_hits        /* adjust number of cache hits */
LN192:  ...
        jbsr    _foo
        cml    #0,_blkmarker+956  /* check if block 239 is in cache */
        jeq     LN193             /* if not then have to invoke simulator */
        cml    #9997,_c_switch    /* check if context switch in 3 cycles */
        jge     LN193             /* if pending then invoke simulator */
        addql   #3,_c_switch      /* adjust context switch information */
        addql   #3,_c_hits        /* adjust number of cache hits */
LN194:  ...
        jne     L405
        movl    #-1,_highmarker   /* nothing to clear since not in loop */

```

Figure 9: Assembly Code with Technique E and Context Switches

The variation of Technique E that assumes there are no context switches is similar to the method illustrated in Figure 7. Not only are the array elements associated with the blocks in the loop checked in the exit blocks, but the array elements associated with the invoked routines are checked in the exit blocks as well. In a function that does not fit in cache, instructions are inserted after calls to a function being measured that does fit in the cache to adjust the number of cache hits and misses. All array elements are checked at the end of the execution of the program in case the program terminated in a routine that fits in cache. Figure 10 illustrates the same loop as in Figure 9 with the assumption that no context switches can occur.

```

L405:  addq1    #5,_blkmarker+952    /* adjust number of hits for block 238 */
      ...
      jbsr    _foo
      addq1   #3,_blkmarker+956    /* adjust number of hits for block 239 */
      ...
      jne     L405
      ...                          /* check blocks 238-239 and blocks in foo

```

Figure 10: Assembly Code with Technique E and No Context Switches

Technique F

To be able to avoid address references being processed by the cache simulator, Technique E required that all the blocks in a loop and the blocks in the routines that can be invoked from the loop fit in the cache at the same time. Technique F relaxes this requirement. All of the blocks in the loop itself still have to fit in the cache at the same time. Also, each individual function invoked directly from the loop and the set of routines the function could in turn invoke have to fit in the cache at the same time. Otherwise there can be conflicts between blocks. For instance, a basic block in a routine invoked from the loop could conflict with a basic block in the loop or a block from a different routine invoked from the same loop. A heuristic was applied that requires that at least one half of the blocks in the loop and the routines invoked cannot conflict. If the heuristic is not satisfied then Technique C is used for the blocks in the loop. Otherwise, if a block in the loop conflicts with blocks in routines invoked from the loop, then instructions are inserted to clear the markers associated with the conflicting blocks in the invoked routines when the loop block is executed. Likewise, when a block in an invoked function conflicts with a loop block or a different routine that is invoked from the same loop, then instructions are inserted to clear the conflicting blocks before the call instruction to the invoked function. Figure 11 shows a loop containing a block that conflicts with a basic block in a routine invoked from the loop.

The variation of Technique F that assumes there are no context switches is accomplished in the following manner. If a block in the loop conflicts with the blocks in routines invoked from the loop, then instructions are inserted to adjust the number of cache hits and misses associated with the conflicting blocks in the invoked routines when the loop block is executed. In a similar fashion, when a block in an invoked function conflicts with a

```

L85:  ...
      cmpl  #0,_blkmarker+2120 /* check if block 530 is in cache */
      jeq   LN141             /* if not then have to invoke simulator */
      cmpl  #9981,_c_switch  /* check if context switch in 19 cycles */
      jge   LN141             /* if pending then invoke simulator */
      addl  #19,_c_switch    /* adjust context switch information */
      addl  #19,_c_hits      /* adjust number of cache hits */
LN142: clrl  _blkmarker+7344  /* block 530 in loop conflicts with block 1836 */
      ...
      clrl  _blkmarker+2120  /* block 1836 in foo conflicts with block 530 */
      jbsr  _foo
      ...
      jne   L85

```

Figure 11: Assembly Code with Technique F and Context Switches

loop block or a different routine that is invoked from the same loop, then instructions are inserted to adjust the number of cache hits and misses associated with the conflicting blocks before the call instruction to the invoked function. Figure 12 shows how the same loop in Figure 11 is instrumented with trace instructions in the variation that assumes there are no context switches.

Technique G

Techniques D-F attempt to find basic blocks that are already resident in the cache because of temporal locality due to loops. There is, however, another situation when temporal locality can result in blocks being resident. This situation occurs when a routine is invoked from more than one location and some blocks in the routine have not been replaced when the second call occurs. Assuming that context switches can occur, Techniques E and F clear the block markers of a function that fits in the cache before the call instruction to that function in a routine that does not fit in cache. At the point when a call to such a function is encountered in Technique G, the compiler attempts to determine if any block markers in the function have already been cleared and the lines associated with those blocks have not been replaced. Any block that is determined to be resident, and thus its associated marker is already cleared, need not have its marker cleared again. The example in Figure 13 has two calls to the same function that fits in the cache. In this example, the instructions following the first call reside in the same line as only one of the blocks in the function being invoked.

```

L85:  ...
      addl   #19, _blkmarker+2120    /* adjust number of hits for block 530 */
      cmpl   #0, _blkmarker+7344    /* check if block 1836 was executed */
      jlt    LN142                  /* if not then skip over measurement code */
      movl   #1836, _startblk        /* first block in sequence */
      pea    #1836                  /* last block in sequence */
      jbsr   _traceblknum           /* call trace routine */
      addq1  #4, a7                  /* adjust stack pointer */
      movl   _blkmarker+7344, d0     /* load remaining hits for block 1836 */
      addl   d0, _c_hits             /* adjust number of cache hits */
      movl   #-7, _blkmarker+7344   /* reset hits for block 1836 */

LN142: ...
      cmpl   #0, _blkmarker+2120    /* check if block 530 was executed */
      jlt    LN143                  /* if not then skip over measurement code */
      movl   #530, _startblk        /* first block in sequence */
      pea    #530                  /* last block in sequence */
      jbsr   _traceblknum           /* call trace routine */
      addq1  #4, a7                  /* adjust stack pointer */
      movl   _blkmarker+2120, d0     /* load remaining hits for block 530 */
      addl   d0, _c_hits             /* adjust number of cache hits */
      movl   #-19, _blkmarker+2120  /* reset hits for block 530 */

LN143: ...
      jbsr   _foo
      ...
      jne    L85

```

Figure 12: Assembly Code with Technique F and No Context Switches

```

      clr1   _blkmarker+528          /* clear blkmarker[132] */
      clr1   _blkmarker+532          /* clear blkmarker[133] */
      clr1   _blkmarker+536          /* clear blkmarker[134] */
      clr1   _blkmarker+540          /* clear blkmarker[135] */
      movl   #132, _lowmarker        /* first block cleared on context switch */
      movl   #135, _highmarker       /* last block cleared on context switch */
      jbsr   _foo
      ...
      clr1   _blkmarker+528          /* clear blkmarker[132] */
      jbsr   _foo
      movl   #-1, _highmarker        /* nothing to clear since not in func */

```

Figure 13: Assembly Code with Technique G and Context Switches

The variation for Technique G that assumes there are no context switches is accomplished differently. After the call to a function that fits in cache, each block marker associated with the called function is examined to update the number of cache hits and misses and the marker is reset to the negated number of references associated with the block. If it is determined that there exists another call to a function that always follows the current call and will result in a block marker being rechecked and the block has not been replaced, then the check after the current call is not inserted. The example in Figure 14 also has two calls to the same function that fits in cache. Only the marker for the conflicting block is checked after the first call.

```

jbsr    _foo
cmlpl  #0,_blkmarker+528    /* check if block 132 was executed */
jlt    LN75                /* if not then skip over measurement code */
movl   #132,_startblk      /* first block in sequence */
pea    #132                /* last block in sequence */
jbsr   _traceblknum        /* call trace routine */
addq1  #4,a7               /* adjust stack pointer */
movl   _blkmarker+528,d0    /* load remaining hits for block 528 */
addl   d0,_c_hits          /* adjust number of cache hits */
movl   #-13,_blkmarker+528 /* reset hits for block 1836 */
LN75:  ...
jbsr   _foo
...
/* now check all four blocks */

```

Figure 14: Assembly Code with Technique G and No Context Switches

RESULTS

The set of test programs used in this experiment and their associated code size are shown in Table 1. The code size for each program does not include the routines from the run-time library since their source code was not available. The techniques discussed in this paper could be used to process assembly or object files. Unfortunately, this would complicate implementation of the techniques since portability would be decreased and the control-flow and data-flow information already available in a compiler would have to be recalculated. A C compiler for the Motorola 68020/68881 was implemented within the *ease* environment^{22,23}, which consists of a compiler generation system called *vpo*²⁴ and measurement tools. The compiler was modified to implement

each of the seven techniques described in the previous section. Cache performance measurements were obtained for each program within the test set using each of the techniques. The measurements obtained for each specific program were not affected by the technique used. Identical results occurred, the exact number of hits and misses, despite periodic simulated context switches and each program requiring at least one million cache references.

Name	Description	Size in Bytes
compact	Huffman Coding Compression	4322
cpp	C Preprocessor	12678
diff	Differences between Files	9166
lex	Lexical Analyzer Generator	26318
sed	Stream Editor	13946
sort	Sort or Merge Files	5500
tbl	Table Formatter	24592
yacc	Parser Generator	22392

Table 1: Test Programs

Periodic context switches were simulated by invalidating the entire cache every 10,000 units of work. A cache hit was assumed to require one work unit while a cache miss was assumed to require ten. The context switch interval and estimated time units required for a hit versus a miss are the same as those used in Smith's cache studies²⁰. Though the experiments in this paper simulated context switching based on estimated cache work to check that identical measurements were obtained with the different techniques, other methods to determine context switch points could also be used.

Table 2 shows the number of times that the cache simulator was invoked for each program using each of the techniques for a 1K byte direct-mapped cache with a 16 byte line size and periodic context switches. The hit ratio is given to indicate the percentage of references that are candidates for not being processed by the cache simulator. Note that the number of calls to the cache simulator using Technique A is the same as the

number of instructions that were executed.* The results using Technique B indicate that there were on average 2.77 (i.e. $1 / .3612$) instructions per basic block being executed. There was also on average 5.29 (i.e. $1 / .1892$) contiguous instructions being processed by the cache simulator using Technique C. Technique D, which requires no inter-procedural analysis, resulted in a substantial improvement over Technique C. This indicates that a large percentage of instructions executed in programs occur in loops with no calls. Techniques E, F, and G appear to be more closely affected by the hit ratio. The results indicate that Technique F has a slight improvement over Technique E. Technique G, however, rarely resulted in fewer references being processed as compared to Technique F. It is interesting to note that using Techniques E, F, and G can occasionally result in more address references being processed by the cache simulator. This situation can occur when the loops that do fit in cache are also in the functions that fit in cache. Since each block in these functions are processed individually, then the blocks outside the loops in these functions would require more references to be processed by the cache simulator since the method used in Technique C is not applied. When no context switches were allowed, the number of calls to the cache simulator decreased slightly with Techniques B-G since the cache simulator never had to process instructions within a block individually.

Program	Hit Ratio	Instructions Executed	Relative to Technique A					
			B	C	D	E	F	G
compact	95.25%	4,699,295	26.51%	13.21%	10.92%	8.03%	8.03%	8.03%
cpp	93.78%	1,322,671	26.53%	15.75%	11.07%	11.01%	11.02%	11.02%
diff	99.54%	3,425,264	29.83%	15.07%	3.22%	1.22%	0.77%	0.77%
lex	99.48%	36,844,880	50.55%	19.47%	2.32%	1.89%	0.94%	0.94%
sed	96.49%	1,643,093	45.86%	26.84%	6.99%	6.92%	5.32%	5.32%
sort	96.91%	1,778,463	34.36%	18.49%	12.24%	12.59%	12.59%	12.59%
tbl	86.59%	2,715,097	30.70%	19.65%	17.28%	16.78%	16.36%	16.24%
yacc	98.63%	23,960,045	44.58%	22.89%	5.17%	4.52%	3.51%	3.51%
average	95.83%	9,548,601	36.12%	18.92%	8.65%	7.87%	7.32%	7.30%

Table 2: Calls to Cache Simulator with a 1K Byte Cache

* The number of executed instructions is slightly less than the total cache references simulated since a Motorola 68020/68881 instruction may span two cache lines.

Table 3 shows the number of times that the cache simulator was invoked with cache sizes ranging from 2K to 16K bytes. Increasing the cache size did not vary the number of references processed by the cache simulator using Technique A since the number of instructions that were executed remained the same. Also, the number of references processed using Techniques B, C, and D varied only slightly as the cache size was increased.* Therefore, only the hit ratio and results from Techniques E, F, and G were presented in Table 3. Slightly varying number of references processed when the cache size was increased for each program using Technique D indicates that loops with no calls in the test set always fit in a 1K byte cache. Unlike Techniques A-D, Techniques E-G improved as the cache size and hit ratios increased. For a program that executes a very large number of instructions, more time may be saved by recompiling the program when the number of sets in the cache configuration being measured is increased. Again the number of calls to the cache simulator decreased only slightly when no context switches were allowed.

Table 4 shows the execution time required using a 1K byte direct-mapped cache with a 16 byte line size and periodic context switches relative to execution without tracing for each program.† Note that execution times of the programs being measured include both the time required for generating the trace and analyzing the references with the cache simulator. The ratio to execution time without tracing for the different programs with each technique varied. For instance, the ratio for Technique A for the `lex` program was over 3.8 times as great as the ratio for `tbl`. The ratio to execution time without tracing is affected by a number of factors which include the average execution time required for the non-tracing instructions executed, the average number of instructions in executed basic blocks, and the percentage of time spent in the library routines which were not measured.‡

* When a context switch could occur in a basic block, the cache simulator processes each instruction within the block individually. Since changing the cache size typically resulted in context switches occurring in basic blocks with a different number of instructions, there was a slight variation in the number of times that the cache simulator was invoked.

† All execution times reported in this paper were obtained by determining the average of ten execution times of each instance of a program.

‡ Smaller ratios to execution times without tracing were reported for a method similar to Technique B in the *trapeds* system¹⁷. This discrepancy was probably due to their choice to simulate more floating-point intensive programs, to not introduce or check for pending context switches, and the use of a cache simulator tuned for specific cache configurations.

Cache Size	Program	Hit Ratio	Relative to Technique A		
			E	F	G
2K bytes	compact	96.04%	8.03%	8.03%	8.03%
	cpp	96.52%	11.03%	11.03%	11.02%
	diff	99.68%	1.22%	0.77%	0.77%
	lex	99.49%	0.94%	0.94%	0.94%
	sed	98.23%	6.92%	3.93%	3.93%
	sort	99.83%	9.91%	5.41%	5.41%
	tbl	93.94%	13.40%	13.42%	13.17%
	yacc	99.28%	2.88%	2.89%	2.89%
average	97.88%	6.79%	5.80%	5.77%	
4K bytes	compact	97.63%	8.03%	2.46%	2.46%
	cpp	97.72%	11.23%	11.07%	11.07%
	diff	99.75%	0.28%	0.28%	0.28%
	lex	99.59%	0.94%	0.94%	0.94%
	sed	98.23%	3.93%	3.93%	3.93%
	sort	99.83%	0.22%	0.23%	0.23%
	tbl	95.69%	13.38%	13.41%	13.16%
	yacc	99.45%	2.86%	2.88%	2.88%
average	98.49%	5.11%	4.40%	4.37%	
8K bytes	compact	99.26%	0.72%	0.72%	0.72%
	cpp	98.06%	10.90%	10.30%	10.30%
	diff	99.75%	0.28%	0.28%	0.28%
	lex	99.60%	0.92%	0.92%	0.92%
	sed	99.11%	0.82%	0.82%	0.82%
	sort	99.82%	0.22%	0.22%	0.22%
	tbl	96.76%	13.41%	13.44%	13.19%
	yacc	99.48%	0.97%	0.99%	0.99%
average	98.98%	3.53%	3.46%	3.43%	
16K bytes	compact	99.26%	0.72%	0.72%	0.72%
	cpp	98.39%	1.36%	1.36%	1.36%
	diff	99.75%	0.28%	0.28%	0.28%
	lex	99.60%	0.44%	0.41%	0.41%
	sed	99.11%	0.82%	0.82%	0.82%
	sort	99.83%	0.22%	0.22%	0.22%
	tbl	97.89%	10.78%	7.05%	6.95%
	yacc	99.48%	0.63%	0.63%	0.63%
average	99.16%	1.91%	1.44%	1.42%	

Table 3: Calls to Cache Simulator with Larger Cache Sizes

Program	Ratio to Execution Time without Tracing						
	A	B	C	D	E	F	G
compact	124.57	51.35	40.23	33.22	25.70	25.74	25.80
cpp	113.05	47.83	40.24	31.55	31.62	31.62	31.67
diff	83.04	35.64	26.51	8.31	5.20	4.50	4.51
lex	188.68	118.70	73.37	15.13	13.68	10.69	10.69
sed	143.18	87.62	72.85	24.94	25.35	22.05	22.05
sort	161.03	83.41	64.75	46.39	47.64	47.72	47.64
tbl	49.65	25.22	22.88	21.01	21.10	20.82	20.35
yacc	124.56	75.12	54.82	17.29	16.46	14.20	14.09
average	123.47	65.61	49.46	24.73	23.34	22.17	22.10

Table 4: Execution Time Overhead with a 1K Byte Cache and Context Switches

Table 5 shows the execution time required using a 1K byte direct-mapped cache with a 16 byte line size and no periodic context switches relative to execution without tracing for each program. The ratios for all the techniques decreased since the trace routine and cache simulator did not have to check for periodic context switches. Techniques D-G also improved since fewer trace instructions were executed inside loops.

Program	Ratio to Execution Time without Tracing						
	A	B	C	D	E	F	G
compact	115.72	45.72	33.28	28.00	20.70	21.10	20.78
cpp	97.57	40.97	33.21	25.13	25.30	25.37	25.34
diff	64.44	27.94	19.51	5.09	2.69	2.03	2.04
lex	139.62	84.40	44.52	7.37	6.02	4.11	4.11
sed	111.16	59.60	46.96	15.61	15.86	13.30	13.76
sort	121.23	55.17	41.93	31.38	31.00	31.65	31.77
tbl	41.83	20.38	17.59	15.67	15.48	15.18	15.27
yacc	94.81	51.46	35.09	10.24	9.30	9.35	9.31
average	98.30	48.21	34.01	17.31	15.80	15.26	15.30

Table 5: Execution Time Overhead with a 1K Byte Cache and No Context Switches

Table 6 shows the execution time overhead with cache sizes ranging from 2K to 16K bytes and periodic context switches. In general, the execution times decrease as the number of references processed by the cache simulator decrease. Thus, as the cache size increases, the execution times for programs using Techniques E, F,

Cache Size	Program	Ratio to Execution Time without Tracing		
		E	F	G
2K bytes	compact	25.58	25.63	25.57
	cpp	29.58	29.60	29.58
	diff	5.12	4.41	4.44
	lex	10.83	10.83	10.75
	sed	16.90	16.90	16.91
	sort	37.72	24.30	24.33
	tbl	15.96	16.32	16.01
	yacc	12.48	12.53	12.52
	average	19.27	17.57	17.51
4K bytes	compact	24.47	11.56	11.51
	cpp	29.65	29.17	29.45
	diff	3.76	3.78	3.81
	lex	11.02	10.93	10.98
	sed	16.98	17.19	17.14
	sort	6.87	8.64	8.55
	tbl	15.42	15.42	15.05
	yacc	12.52	12.83	12.68
	average	15.09	13.69	13.65
8K bytes	compact	6.67	6.67	6.67
	cpp	29.14	27.93	27.95
	diff	4.08	4.08	4.08
	lex	11.38	11.47	11.38
	sed	9.49	9.59	9.61
	sort	7.63	7.64	7.65
	tbl	15.22	15.21	15.29
	yacc	9.08	9.11	9.12
	average	11.59	11.46	11.47
16K bytes	compact	7.44	7.44	7.45
	cpp	8.99	9.02	8.97
	diff	4.73	4.73	4.71
	lex	11.91	11.74	11.74
	sed	10.95	10.97	10.97
	sort	8.58	8.60	8.58
	tbl	13.29	10.14	10.02
	yacc	8.96	8.97	8.92
	average	9.36	8.95	8.92

Table 6: Execution Time Overhead with Larger Cache Sizes and Context Switches

and G decrease. The execution times for Techniques A-D varied only slightly since the number of references processed by the cache simulator only changed slightly with different cache sizes simulated. Therefore, only the execution time ratios for Techniques E, F, and G are presented in Table 6.

Table 7 shows the execution time overhead for Techniques E, F, and G with cache sizes ranging from 2K to 16K bytes and no periodic context switches. As the cache size increases, the reduction in execution time resulted in greater improvements as compared to the execution times with context switches. These greater improvements were due to the high percentage of blocks that fit within loops and fewer trace instructions being executed with the variations that do not allow periodic context switches.

The tracing overhead is dependent on the performance of the cache simulator. Less tracing overhead would be required if cache simulators were used that were tuned for a particular cache configuration. This scheme, however, would be less flexible since a program would have to be relinked each time the cache configuration was changed.

FUTURE WORK

The techniques described in this paper could be extended to have the ability to evaluate instruction cache performance from programs executing concurrently. The method used by the techniques to simulate periodic context switches had the advantage that it disregarded the time required to execute the instructions that were inserted to obtain the measurements. This advantage is important since the measurement overhead not only dominates the total execution time, but also can vary during different portions of the execution because calls to the trace routine are not uniformly distributed. Also, context switches were allowed to occur between instructions within a basic block, which is more realistic than methods that only allow context switches on basic block boundaries. The effect of each context switch was simulated by invalidating the entire instruction cache. Unfortunately, most current systems do not cold start the cache when switching to a new process. The techniques described in this paper could be extended to measure cache performance for multiple concurrently executing programs in the following manner. First, the data used to represent the state of the cache in the cache

Cache Size	Program	Ratio to Execution Time without Tracing		
		E	F	G
2K bytes	compact	20.19	20.25	20.20
	cpp	24.27	24.48	24.84
	diff	2.55	1.95	1.93
	lex	3.99	3.94	3.97
	sed	9.55	9.65	9.82
	sort	24.34	13.81	14.01
	tbl	12.12	12.24	12.24
	yacc	5.91	5.91	5.91
	average	12.87	11.53	11.61
4K bytes	compact	20.00	6.07	5.82
	cpp	25.06	24.91	24.82
	diff	1.26	1.26	1.24
	lex	3.81	3.79	3.85
	sed	8.90	8.67	8.79
	sort	1.46	4.79	4.84
	tbl	11.86	11.94	11.76
	yacc	5.88	5.90	5.93
	average	9.78	8.42	8.38
8K bytes	compact	1.51	1.53	1.53
	cpp	25.17	24.38	24.46
	diff	1.25	1.25	1.26
	lex	3.86	3.91	3.89
	sed	1.68	1.63	1.67
	sort	1.76	1.78	1.74
	tbl	10.63	10.68	10.65
	yacc	2.54	2.57	2.56
	average	6.05	5.97	5.97
16K bytes	compact	1.53	1.52	1.53
	cpp	1.45	1.45	1.44
	diff	1.24	1.24	1.24
	lex	2.09	1.77	1.77
	sed	1.68	1.69	1.69
	sort	1.69	1.73	1.67
	tbl	8.90	6.38	6.20
	yacc	1.78	1.77	1.77
	average	2.55	2.19	2.16

Table 7: Execution Time Overhead with Larger Cache Sizes and No Context Switches

simulator would be mapped to a shared global data area. Rather than invalidate the entire cache at a periodic context switch point, the cache simulator will suspend itself until it receives a signal indicating that the current process can continue. At this point control will be passed to the next program in a round robin fashion. At the desired termination point, instruction cache measurements could be generated for each individual program and/or for the entire set of programs.

Another area of future work would be to determine if the techniques discussed in this paper could be applied to machines that fetch and issue multiple instructions at the same time. (V)LIW machines would not be difficult to handle since the compiler packages the program into sets of instructions that are always fetched and issued together. The challenge would be determining the sets of instructions that are actually fetched and issued simultaneously for superscalar machines since these sets would depend on the paths taken through a program. The number of instructions issued together would also depend not only structural hazards (available functional units), but also the alignment requirements of a fetch block²⁵.

CONCLUSIONS

The techniques presented in this paper have been shown to significantly reduce the time required for instruction cache performance evaluations as compared to more traditional approaches. This improvement occurred despite no special requirements to implement the techniques* and without any loss of accuracy. Technique D is particularly attractive since with no interprocedural analysis required it is simple to implement and still results in a significant improvement. Though only the number of instruction references to be processed can be reduced, the techniques can also be used when evaluating split instruction and data caches. There still should be a measurable improvement in this situation since typically the majority of address references being processed are instructions². The effective evaluation of large second-level caches may require billions of references to be traced. When positioned behind a split first-level cache, the techniques presented in this paper would be very useful.

* Some approaches have dedicated a set of registers to be used exclusively for tracing and/or require special operating system support^{2,7,9,11}.

REFERENCES

1. A. J. Smith, 'Two Methods for the Efficient Analysis of Memory Address Trace Data', *IEEE Transactions on Software Engineering*, **3**, 94-101 (January 1977).
2. A. Borg, R. E. Kessler and D. W. Wall, 'Generation and Analysis of Very Long Address Traces', *Proceedings of the 17th Annual International Symposium on Computer Architecture*, Seattle, Washington, 270-279 (May 1990).
3. W. W. Hwu and P. P. Chang, 'Achieving High Instruction Cache Performance with an Optimizing Compiler', *Proceedings of the 16th Annual Symposium on Computer Architecture*, Jerusalem, Israel, 242-250 (May 1989).
4. S. McFarling, 'Program Optimization for Instruction Caches', *Proceedings of the Third International Conference on Architectural Support for Programming Languages and Operating Systems*, Boston, Ma, 183-191 (April 1989).
5. K. Pettis and R. Hansen, 'Profile Guided Code Positioning', *Proceedings of the SIGPLAN Notices '90 Conference on Programming Language Design and Implementation*, White Plains, New York, 16-27 (June 1990).
6. B. L. Peuto and L. J. Shustek, 'An Instruction Timing Model of CPU Performance', *Proceedings of the 4th Annual Symposium on Computer Architecture*, Silver Spring, Maryland, 165-178 (March 1977).
7. C. A. Wiecek, 'A Case Study of VAX-11 Instruction Set Usage for Compiler Execution', *Proceedings of the Symposium on Architectural Support for Programming Languages and Operating Systems*, Palo Alto, California, 177-184 (March, 1982).
8. M. Huguet, T. Lang and Y. Tamir, 'A Block-and-Actions Generator as an Alternative to a Simulator for Collecting Architecture Measurements', *Proceedings of the SIGPLAN Notices '87 Symposium on Interpreters and Interpretive Techniques*, St. Paul, Minnesota, 14-25 (June 1987).
9. A. Agarwal, R. L. Sites and M. Horowitz, 'ATUM: A New Technique for Capturing Address Traces Using Microcode', *Proceedings of the 13th Annual Symposium on Computer Architecture*, Tokyo, Japan, 119-127 (June 1986).
10. D. W. Clark and H. M. Levy, 'Measurement and Analysis of Instruction Use in the VAX-11/780', *Proceedings of the 9th Annual Symposium on Computer Architecture*, Austin, Texas, 9-17 (April 1982).
11. S. J. Eggers, D. R. Keppel, E. J. Kolding and H. M. Levy, 'Techniques for Efficient Inline Tracing on a Shared-Memory Multiprocessor', *Proceedings SIGMETRICS '90 Conference on Measurement and Modeling of Computer Systems*, Boulder, CO, 37-47 (May 1990).
12. T. Ball and J. R. Larus, 'Optimally Profiling and Tracing Programs', *Conference Record of the Nineteenth Annual Symposium on Principles of Programming Languages*, Albuquerque, NM, 59-70 (January 1992).
13. T. R. Puzak, *Analysis of Cache Replacement Algorithms*, PhD Dissertation, University of Massachusetts, Amherst, MA, February 1985.
14. W. Wang and J. Baer, 'Efficient Trace-Driven Simulation Methods for Cache Performance Analysis', *Proceedings SIGMETRICS '90 Conference on Measurement and Modeling of Computer Systems*, Boulder, CO, 27-36 (May 1990).
15. M. D. Hill and A. J. Smith, 'Evaluating Associativity in CPU Caches', *IEEE Transactions on Computers*, **38**, 1612-1630 (December 1989).
16. S. Laha, J. H. Patel and R. K. Iyer, 'Accurate Low-Cost Methods for Performance Evaluation of Cache Memory Systems', *IEEE Transactions on Computers*, **37**, 1325-1336 (November 1988).

17. C. Stunkel and W. Fuchs, 'TRAPEDS: Producing Traces for Multicomputers Via Execution Driven Simulation', *Proceedings of the International Conference on Measurement and Modeling of Computer Systems*, , 70-78 (May 1989).
18. C. L. Mitchell and M. J. Flynn, 'A Workbench for Computer Architects', *IEEE Design & Test of Computers*, **5**, 19-29 (February 1988).
19. D. B. Whalley, 'Fast Instruction Cache Performance Evaluation Using Compile-Time Analysis', *Proceedings SIGMETRICS and PERFORMANCE '92 Conference on Measurement and Modeling of Computer Systems*, Newport, RI, 13-22 (June 1992).
20. A. J. Smith, 'Cache Memories', *Computing Surveys*, **14**, 473-530 (September 1982).
21. J. Hennessy and D. Patterson, *Computer Architecture: A Quantitative Approach*, Morgan Kaufmann, San Mateo, CA, 1990.
22. J. W. Davidson and D. B. Whalley, 'Ease: An Environment for Architecture Study and Experimentation', *Proceedings SIGMETRICS '90 Conference on Measurement and Modeling of Computer Systems*, Boulder, CO, 259-260 (May 1990).
23. J. W. Davidson and D. B. Whalley, 'A Design Environment for Addressing Architecture and Compiler Interactions', *Microprocessors and Microsystems*, **15**, 459-472 (November 1991).
24. M. E. Benitez and J. W. Davidson, 'A Portable Global Optimizer and Linker', *Proceedings of the SIGPLAN Notices '88 Symposium on Programming Language Design and Implementation*, Atlanta, GA, 329-338 (June 1988).
25. M. D. Smith, M. Johnson and M. A. Horowitz, 'Limits on Multiple Instruction Issue', *Proceedings of the Third International Conference on Architectural Support for Programming Languages and Operating Systems*, Boston, Ma, 290-302 (April 1989).