

Cache Performance of the SPEC92 Benchmark Suite

Jeffrey D. Gee

Sun Microsystems Inc.

Mark D. Hill

Dionisios N. Pnevmatikatos

Computer Sciences Department
University of Wisconsin-Madison
Madison, WI 53706

Alan Jay Smith

Computer Science Division
Dept. of Electrical Engineering
and Computer Sciences
University of California Berkeley
Berkeley, CA 94720

This paper will appear in IEEE Micro

Abstract

The SPEC92 benchmark suite consists of twenty public-domain, non-trivial programs that are widely used to measure the performance of computer systems, particularly those in the Unix workstation market. These benchmarks were expressly chosen to represent real-world applications and were intended to be large enough to stress the computational and memory system resources of current-generation machines. The extent to which the *SPECmarks* (the figures of merit obtained from running the SPEC benchmarks under certain specified conditions) accurately represents performance with live real workloads is not well established; in particular, there is some question whether the memory referencing behavior (cache performance) is appropriate.

In this paper, we present measurements of miss ratios for the entire set of SPEC92 benchmarks for a variety of CPU cache configurations; this study extends earlier work that measured only the performance of the integer SPEC89 benchmarks. We find that instruction cache miss ratios are generally very low, and that data cache miss ratios for the integer benchmarks are also quite low. Data cache miss ratios for the floating point benchmarks are more in line with published measurements for real (i.e. non-benchmark, non-synthetic) workloads. We believe that the discrepancy between the SPEC benchmark miss ratios and those observed elsewhere is partially due to the fact that the SPEC benchmarks are all almost exclusively user state CPU benchmarks run until completion as the single active user process. We therefore believe that SPECmark performance levels may not reflect system performance when there is multiprogramming, time sharing and/or significant operating systems activity.

KEYWORDS: Cache performance, Memory system design, Multiprogramming, Program behavior, SPEC benchmarks, Trace-driven simulation

1. Introduction

The 1992 SPEC benchmarks [SPEC91] are a selection of non-trivial programs chosen to standardize benchmarking. SPEC (System Performance Evaluation Consortium) assembled this suite to provide a standard set of realistic benchmarks for inter-system comparisons; see [Pric89, Hinn88] for a discussion of the many problems with the benchmarking situation prior to SPEC. Several factors, including strong industrial support for SPEC, the realistic nature of the benchmarks, and acceptable code portability have led to the wide use of these programs for benchmarking purposes. To improve the verification and reproducibility of results, SPEC benchmark results must include a description of any source code modifications, compiler and operating system release numbers, machine characteristics, and most other factors that can affect the reported results. The SPEC benchmarks have become so important as a measure of CPU performance that some system developers are parameterizing their designs to maximize SPEC benchmark performance, even when this might lead to lower performance on other, perhaps more realistic, workloads. Similarly, compiler writers have been concentrating on producing good code for the frequently executed inner loops of some of the SPEC benchmarks. Recent very high benchmark results for the SPEC89 program *matrix300* demonstrated the success of their efforts and forced the SPEC consortium not to include *matrix300* in the SPEC92 benchmark release.

SPEC Benchmark Suite		
Program	Language	Description
alvinn	C	Robotics neural network training
compress	C	Adaptive Lempel-Ziv compression
doduc	Fortran	Thermohydraulic simulation of a nuclear reactor
ear	C	Human ear simulation
eqntott	C	Builds truth table from a boolean expression
espresso	C	Boolean function minimization
fpppp	Fortran	Two electron integral derivative
gcc	C	GNU C compiler compiling pre-processed source files
hydro2d	Fortran	Galactical jet computation
mdljdp2	Fortran	Molecular dynamics (double precision)
mdljsp2	Fortran	Molecular dynamics (single precision)
nasa7	Fortran	Seven floating-point synthetic kernels
ora	Fortran	Ray tracing
sc	C	Spreadsheet calculator
spice	Fortran	Analog circuit simulator
su2cor	Fortran	Quantum physics mass computation
swm256	Fortran	Shallow water equation solver
tomcatv	Fortran	Mesh generation program
wave5	Fortran	Maxwell's equation solver
xlisp	C	Lisp interpreter solving the nine queens problem

Table 1: SPEC Benchmark Applications

The SPEC92 benchmark suite consists of six integer-intensive C programs (*compress*, *eqntott*, *espresso*, *gcc*, *sc*, and *xlisp*) and fourteen floating-point intensive programs (*alvinn*, *doduc*, *ear*, *fpppp*, *hydro2d*, *mdljdp2*, *mdljsp2*, *nasa7*, *ora*, *spice*, *su2cor*, *swm256*, *tomcatv*, and *wave5*). The SPEC benchmarking procedure is to run each program to completion on the target system, with only one user process active, and then take the ratio of that run time to the run time of the same program on a DEC VAX 11/780, as measured originally at the start of the SPEC effort. The geometric mean of those ratios over the integer and floating point intensive programs, yields the "*SPECint92*" and "*SPECfp92*" respectively, which are the figures of merit. Table 1 lists and gives a short description of each benchmark.

As noted above, considerable effort is being expended on creating computer systems (hardware and software) to optimize SPEC benchmark results. Two questions therefore arise: (a) In what ways should the system be designed to perform well on the SPEC benchmark suite? (b) Is this a good idea?

One important aspect of CPU performance, and probably the most important of the architectural aspects (as opposed to technology parameters, such as circuit speed) is the performance of the memory hierarchy. We note that SPEC benchmark results are quite sensitive to cache size, as may be seen by comparing the various published measurements of systems, with varying caches sizes, based on the Motorola 88000. (*reference?*) In terms of the SPEC benchmarks, the two questions above become: (a) What miss ratios can be expected when running the SPEC benchmarks on a machine with a cache of a given design? (b) Are these miss ratios comparable to those for "typical" user workloads, for some definition of typical?

In this paper we present measurements of the cache miss ratios of the entire SPEC92 benchmark suite and comment on their potential use in the design of caches and memory hierarchies. We compare the SPEC cache miss ratios to design target miss ratios [Smit87], miss ratios measured using hardware monitors at Amdahl [Smit82] and on DEC VAX-series machines [Clar83,88], miss ratios observed from very long address traces [Borg90], and other miss ratios that include operating system and multiprogramming behavior. We note that miss ratios for multiprogrammed workloads with significant operating system activity are known to be high [Agar88,Ande91]. We find that the miss ratios for the SPEC benchmarks are generally lower than should be expected from multiprogrammed workloads.

2. SPEC Cache Performance

2.1. Methodology

We compiled and ran the SPEC programs on DECstation which contain the MIPS R2000 and R3000 microprocessor, running version 4.1 of the DEC Ultrix operating system. We used version 2.0 of the C compiler and version 2.1 of the Fortran compiler with optimization level according to the SPEC Makefiles. We then used the MIPS *pixie* [DEC91] tool to generate address traces to feed directly to the *tycho* [Hill] cache simulator. *Pixie* modifies the compiled code to generate a trace record for each load, store and basic block entry; trace records for all instruction fetches are then constructed from the basic block records. *Tycho* uses algorithms that, for a given block size, simulate all cache sizes and associativities in a single pass through an address trace [Hill87]. Note that since our traces are derived from the MIPS architecture, different results will be obtained for other CPUs and other compilers.

We varied cache size from 1 Kbyte to 1 Mbyte, set size from one (direct-mapped) to eight, and block size from 16 to 256 bytes. All caches used the LRU replacement algorithm and the lowest order available address bits to select the set. We simulated instruction, data, and unified caches, without any periodic cache flushing, as the SPEC benchmarks are typically run in a uniprogrammed environment. Miss ratios represent the complete execution of a benchmark and include start-up as well as steady-state effects. The use of *pixie* to generate address traces allows simulation of only user, and not system references, and our data is for user code only. Table 2 shows the user and system times for an execution of each of the benchmarks when run on a DECstation 5000/240 MIPS processor based workstation. The system time accounts for 1.5% of the total run time for the benchmark suite, and the linear average of the percentage of system time for each benchmark is 2.5%. The fraction of system time is sufficiently low that we believe that user state only measurements of cache miss ratios are a very accurate approximation of the miss ratios when including both user and system state memory accesses.

Table 2 also lists the number of instruction, data, and total user memory references made by each program. The SPEC92 release specifies that *compress* is run twenty times with the same input and *gcc* is run four times with the same input. The number of references reported here corresponds to one of these runs. Note that the trace reflects a 4-byte memory interface; the trace would be different for a different memory interface width. Note also that the trace includes only actual program loads, stores and instruction fetches; it does not include the extra memory activity such as instruction prefetch that would occur on most machines [Clar83]. For analysis of some of the benchmark programs and their execution behavior, see [Saav92a,b].

To increase our confidence in our results we compared them with two other studies that ran the SPEC benchmarks on a MIPS R2000 microprocessor. Pnevmatikatos and Hill [Pnev90] presented cache miss ratios for the four integer SPEC89 benchmarks (eqntott, espresso, gcc and xlist). They used a different compiler (gcc) and a tracing methodology that excludes library references. Nevertheless, most miss ratio differences are less than 0.01. In few cases, however, a seemingly small miss ratio difference translates into a substantial relative change. We are inclined to place the most confidence in the results presented here, since this analysis has used much more mature and sophisticated compilers, but the comparison demonstrates that cache miss ratios, instruction

Program	Instruction	Data	Total	UserTime	SysTime
Compress	87,493,425	23,079,252	110,572,677	5.5	0.6
Eqntott	1,241,913,236	215,772,134	1,457,685,370	41.1	0.5
Espresso	2,899,136,916	642,332,818	3,541,469,734	81.3	0.7
Gcc	1,262,492,069	398,952,157	1,661,444,226	55.8	5.1
Sc	3,872,103,933	1,180,851,591	5,052,955,524	129.5	18.0
Xlisp	6,257,593,610	2,313,405,716	8,570,999,326	205.6	10.6
Alvinn	6,563,724,007	1,881,599,643	8,445,323,650	250.5	3.1
Doduc	1,619,374,300	583,667,566	2,203,041,866	75.4	0.7
Ear	16,808,813,786	3,973,854,638	20,782,668,424	562.8	5.1
Fpppp	7,420,830,444	4,511,281,096	11,932,111,540	443.9	3.1
Hydro2d	8,398,925,572	3,227,708,646	11,626,634,218	515.6	4.7
Mdljdp2	4,767,422,316	2,063,639,524	6,831,061,840	202.0	0.8
Mdljsp2	3,980,213,579	1,100,297,514	5,080,511,093	159.7	0.6
Nasa7	9,195,719,149	4,720,515,938	13,916,235,087	805.7	14.1
Ora	6,461,088,985	1,606,039,952	8,067,128,937	302.2	1.2
Spice	28,696,843,509	8,288,246,353	36,985,089,862	1496.2	21.2
Su2cor	5,932,445,133	2,579,035,906	8,511,481,039	387.2	2.0
Swm256	11,551,490,879	3,716,782,216	15,268,273,095	662.3	3.5
Tomcatv	1,872,460,468	913,221,318	2,785,681,786	132.3	3.5
Wave5	3,704,008,705	1,224,717,917	4,928,726,622	165.4	1.2
Total	132,594,094,021	45,165,001,895	177,759,095,916	6680.0	100.3
Arithmetic mean				334.0	5.0
Geometric mean				191.5	2.5

Table 2: Program Reference Counts and Execution Time in seconds

counts, and related measures are, as might be expected, sensitive to the compiler used. We must thus caution readers that *your actual mileage may vary*. Cmelik et al. [Cmel91] give instruction counts for the SPEC89 benchmarks. With one exception, Spice, their counts are close to ours. We cannot explain the difference for Spice, although simulation runs at both Berkeley and Madison yielded consistent results.

Simulating these caches required 200 to 400 microseconds of CPU time per memory reference in each trace. Assuming an average 300 microseconds per memory reference, simulating all twenty SPEC benchmarks requires some 980 days or *nearly 40 months* of CPU time. Including false starts, simulation errors, and operating system bugs, we used *three to four years of machine time* to compute our results; this type of measurement would not have been possible if it had been necessary to pay for CPU time on a timeshared machine. (Workstations aren't free, but they are a lot cheaper than the same number of cycles on a timeshared machine.) With seven machines available for running simulations at Berkeley and Madison, we were able to generate these results in less than seven months of calendar time.

2.2. Results

In our simulations we varied the block (line) size from 16 to 256 bytes, the cache size from 4Kbytes to 1 Mbyte, and the set-associativity from 1 (direct mapping) to 8. for instruction, data and unified caches. The complete set of results is not included for brevity reasons, but an electronic copy is available via anonymous ftp[‡]. In this section, we comment on some of that data; in section 3, we present and discuss the averages over the benchmark programs.

We first examine instruction cache miss ratios for the different programs. For *alvinn*, *compress*, *ear*, *eqntott*, *hyrdo2d*, *mdljdp2*, *mdljsp2*, *nasa7*, *ora*, *swm256*, and *tomcatv*, instruction cache miss ratios are very low, generally less than 0.0001 for caches as small as a few kilobytes. These programs spend much of their execution time in a few small routines; the SPEC89 program *matrix300*, for example, spends about 99% of its execution time in one small basic block in the code [Saav90,Saav92a,b]. Miss ratios for *sc*, *espresso*, *su2cor*, *xlisp*, *spice* and *wave* are only slightly larger, as miss ratios again fall below 0.0001 for cache sizes as small as 16 or 32 Kbytes. Instruction cache miss ratios are largest for *doduc*, *gcc*, and *fpppp*, yet are well below half a percent for caches as small as 64 or 128 Kbytes. None of the SPEC benchmarks makes significant use of more than 128 Kbytes of instruction cache.

Miss ratios for data caches are larger, especially for several of the floating-point Fortran benchmarks, but for the most part are quite low as cache size approaches one megabyte. Miss ratios for *ora*, *fpppp*, *xlisp* and *doduc* are the lowest among the SPEC suite, dropping below one percent for caches as small as 16 or 32 Kbytes, and falling below 0.0001 for a 64 Kbyte cache. Results for *ear*, *mdljdp2*, and *espresso* are also low, especially when the set size is greater than one, and for somewhat larger for direct-mapped caches. Among the integer programs, *compress*, *eqntott* and *gcc* exercise fairly large data caches; miss ratios remain above one percent until cache size reaches 512 Kbytes.

The floating-point programs *nasa7*, *spice*, *su2cor*, *swm256*, *tomcatv*, and *wave5* exhibit the largest data cache miss ratios. Miss ratios for *su2cor*, *nasa7*, *spice*, and *wave5* are several percent until the cache size reaches one megabyte, causing miss rates to fall below one percent. *swm256* and *tomcatv* require extremely large caches when the cache block size is small. Data cache miss ratios are *over 12 percent and 6 percent* respectively for a 1 Mbyte cache at a 16-byte block size. Each successive doubling of block size at 1Mbyte reduces data cache miss ratios by almost half, and miss ratios do become less than one percent for a 128 byte block for *tomcatv* and a 256 byte block for

[‡] A machine readable copy of the complete set of tables can be available via anonymous ftp from reggiano.cs.wisc.edu:

```
ftp reggiano.cs.wisc.edu (or: ftp 128.105.8.27 )
reply to login: anonymous
reply to passwd: type any non-null string here
cd SPEC92
get README
get fullmisrratios.ascii
get fullmisrratios.postscript.Z
bye
```

swm256.

Unified (data and instruction) cache miss ratios usually fall between instruction and data cache miss ratios, as the strong locality in instruction references offsets the weaker locality in data references. We do observe several instances where unified cache miss rates are *higher* than corresponding data cache miss rates (*doduc*, *fpppp*, *ora*, *xlisp*). This behavior occurs mainly at larger cache sizes coupled with low associativities, and where separate instruction and data cache miss ratios have fallen to nearly zero. The low associativity causes instruction and data references to conflict for cache sets, while such conflicts do not occur in separate instruction and data caches. Note that a split direct-mapped instruction/data cache pair is more like a 2-way set-associative unified cache than a direct-mapped unified cache.

It is worth noting that there are a few anomalies in the data with respect to the effect of associativity on miss ratio. Generally, miss ratios decrease with increased degrees of set associativity, since the probability of mapping conflicts decreases [Hill89]. It is possible, however, that miss ratios can increase with increasing associativity if certain reference patterns are present in the memory reference string; we note just that effect at one or more data points for the *fpppp*, *spice*, *tomcatv*, and *doduc* miss ratios.

3. Evaluation

In this section we compare the SPEC miss ratios with miss ratios from previous studies and discuss whether the SPEC applications make suitable cache benchmarks. We first describe the other studies.

- (a) Smith [Smit82] includes several measurements taken with a hardware monitor at Amdahl Corporation on various models of the Amdahl 470V machines, running a standard internal benchmark containing supervisor, commercial and scientific code. Results showed that supervisor state miss ratios were much higher than problem state miss ratios, and that the miss ratio for each of user and supervisor state could be approximated by equations of the form $m = a * k^b$, where a and b are constants and k is the cache size in kilobytes.
- (b) Two studies [Clar83,88] provide cache miss ratios taken via hardware measurement from VAX 11/780 and VAX 8800 computers. The 11/780 has an 8 Kbyte, write-through unified cache with an 8-byte block size and a set size of two. The 8800 has a 64 Kbyte, write-through, direct-mapped unified cache with a 64-byte block size. In both cases, these were timeshared workloads, measured at DEC in an engineering environment.
- (c) Smith [Smit85] introduced the *design target miss ratios* (*dtmrs*) to represent typical levels of performance, averaged over a wide class of workloads, ranging from workstations to timeshared mainframes. (In practice, miss ratios for workstations would probably be lower, and for large timeshared mainframes would probably be higher.) He synthesized them from real (hardware monitor) measurements that existed in the literature and a large number of trace-driven simulation results. The initial *dtmrs* for 16-byte line size, fully-associative caches [Smit85] were later extended to other line sizes [Smit87] and to set-associative caches [Hill87,89].

- (d) Agarwal, et al. [Agar88] presented miss ratios that include the effects of operating system references and multiprogramming by using microcode to capture address traces from multi-tasked machines. These effects can more than double miss rates from those measured in a uniprogrammed, user-only environment. They used a varied set of 20 applications programs.
- (e) Borg, et al. [Borg90] generated miss ratios for very long address traces using tools similar to our own; those traces were over twelve billion memory references long. The traces were used to evaluate the performance of a variety of caches. They used three individual traces and another which was a multiprogramming workload consisting of several jobs.

It is important to note that although some of these studies are rather old, we have been unable to find newer or better data. There are many other studies using traces, but we believe that those other workloads are no more representative. With regard to the discussions below, we believe that were any of the above real measurements to be repeated today, the programs and memories would be larger, and the miss ratios (for a given size cache) would be higher.

Figures 1 through 3 show average SPEC miss ratios for instruction, data, and unified caches, with 32-byte lines and 2-way set-associativity, computed separately for the integer and floating-point benchmarks. We also list in Tables 3 through 5 average miss ratios for the integer, floating-point, and complete SPEC92 suite across the entire range of simulation parameters. These averages represent the unweighted arithmetic mean of individual program miss ratios[†]. The unweighted arithmetic mean of the program miss ratios gives the miss ratio of a workload where each program runs for the same number of references. In Figures 1 and 2, average miss rates are plotted against the design target miss ratios (labeled *dtmr*) and primary cache miss ratios from [Borg90] for a multiprogrammed workload (labeled *borg*). Unfortunately, miss ratios from the other studies are not available for separate instruction and data caches, but are plotted against SPEC unified cache results in Figure 3. Previous results based on different block sizes (VAX 11/780, VAX 8800, Agarwal, et al.) or different associativities (VAX 8800, Borg et al.) have been adjusted for these parameters using ratios of miss ratios from prior studies [Hill89, Smit87].

A look at Figure 1 suggests that instruction cache miss ratios for the SPEC benchmarks are unusually low, as they are as low as one-fourth of the design target miss ratios and one-half of Borg's miss ratios.

In Figure 2, we see that data cache miss ratios for the SPEC integer and floating point benchmarks bracket the dtmrs for small cache sizes and are close for the larger sizes for which the dtmrs are defined; all of them are above the [Borg90] measurements. Both sets of SPEC benchmarks approach zero miss ratio for moderately large caches; we would not expect the miss ratios in a timeshared system to approach zero until the cache were as large as main memory because of misses due to task switching (i.e. cold start).

[†] The average miss ratio is calculated using the formula $\frac{1}{n} \sum_{i=1}^n \frac{M_i}{R_i}$ where n is the total number of programs, M_i is the number of misses for program i and R_i is the number of references for program i .

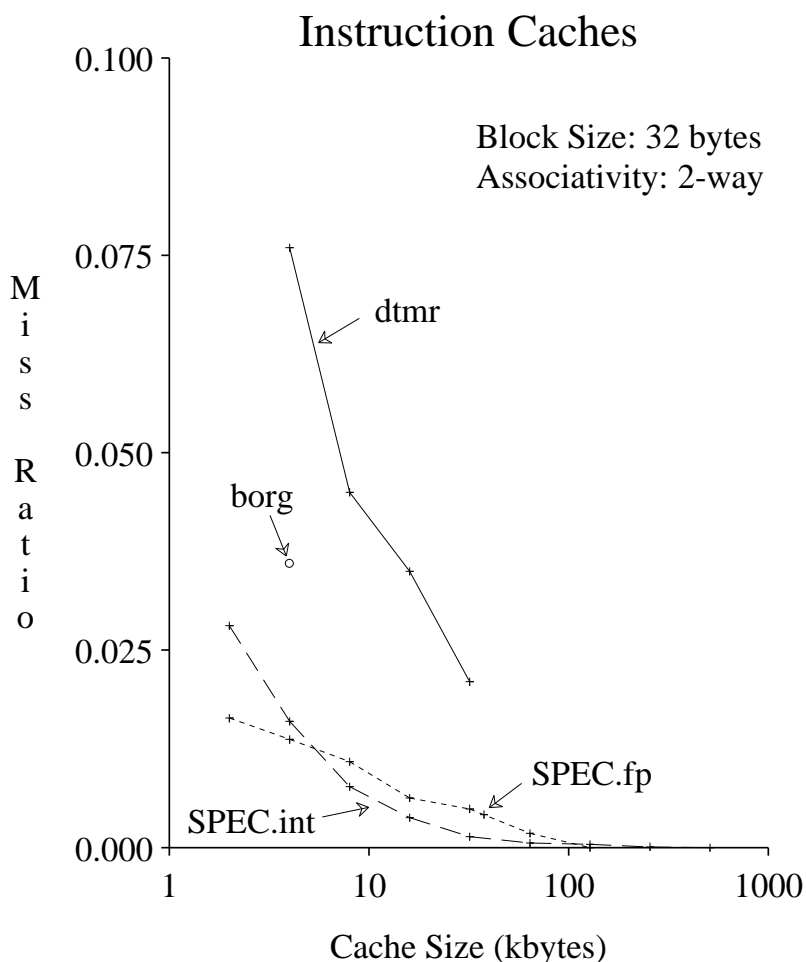


Figure 1: Instruction cache miss ratios

Were the cache the same size as main memory, “misses” would appear as I/O activity, but would still occur.

Figure 3 contains unified cache measurements from the various other studies in addition to SPEC and design target miss ratios. These include: Amdahl 470 supervisor and user state miss ratios (plots labeled *470.sup* and *470.user*), VAX 11/780 and VAX 8800 miss ratios (plots labeled *VAX.780* and *VAX.8800*), and miss ratios from [Agar88] for a multiprogramming level of 3 (plots labeled *agarwal.mul3*). (We plot the Amdahl data from the fitted curve in [Smit82]; the original data points are not available.) We note that the VAX8800 data was collected from a very heavily used timeshared system. The Amdahl 470 supervisor data was collected from the execution of a standard internal Amdahl commercial workload. For both the VAX8800 and Amdahl data, the level of supervisor activity was quite high. Following in decreasing order of miss ratio are the dtmrs and Agarwal’s multiprogrammed miss ratios. SPEC floating point, VAX 11/780 and Amdahl 470 user state miss ratios follow, and the SPEC integer miss rates are smallest by a wide margin.

All of the data in the literature (see e.g. [Smit82], [Ande91], [Agar88]) suggests that operating systems activity significantly increases miss ratios. First, operating

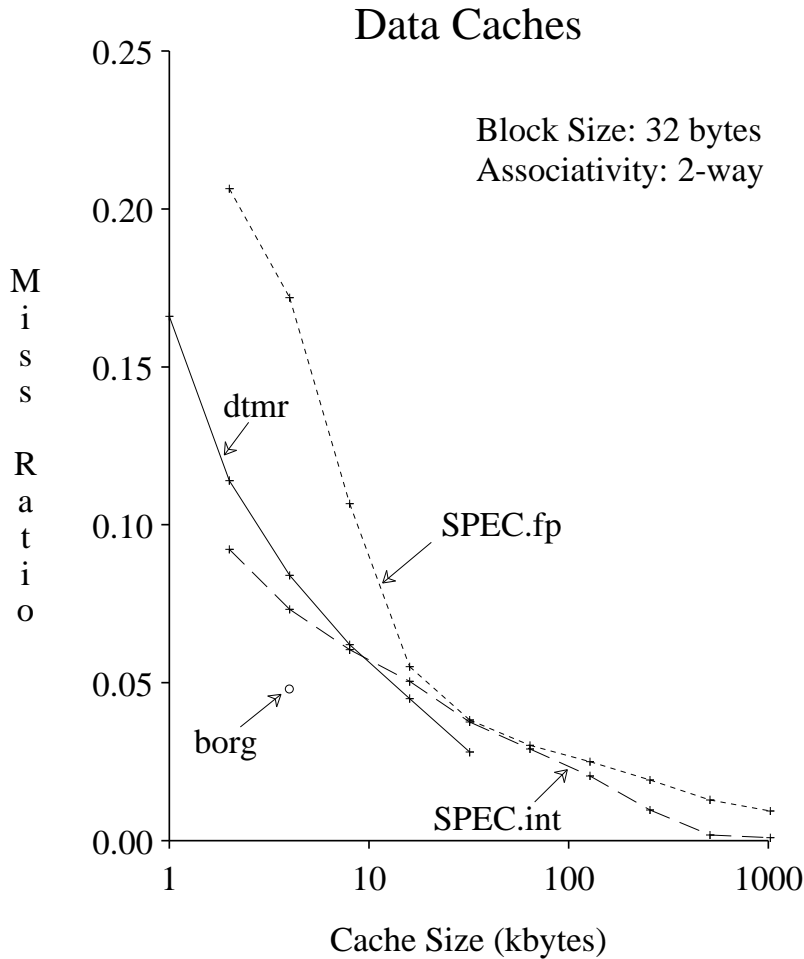


Figure 2: Data cache miss ratios

systems code tends to loop less than user code, and so instruction miss ratios are high. Second, operating systems routines are usually called into the cache by an exception, interrupt or trap, then run for a short time, and finally are replaced from the cache before they run again; they effectively always face a "cold start" situation. Sanguinetti observes [Sang84] that for the Amdahl 580, routines must execute over 600 times per second to stay cache resident. Third, operating system activity is associated with timesharing and high levels of multiprogramming; frequent task switching means that programs are constantly experiencing cold start. As illustrated by Figure 3, miss ratios for the SPEC benchmarks are considerably below those for any workload with significant OS activity, and as noted earlier and as shown in Table 2, the SPEC benchmarks actually contain very little operating system activity. Similar differences in cache performance between compute bound and multiprogrammed environments are reported in [Mogu91]. The SPEC floating-point benchmark miss ratios are quite close to the dtmrs, the data from [Agar88], and the VAX 11/780 measurements, and for large cache sizes are also very close to the Amdahl 470 user program miss ratios. The SPEC integer benchmark miss ratios are the lowest.

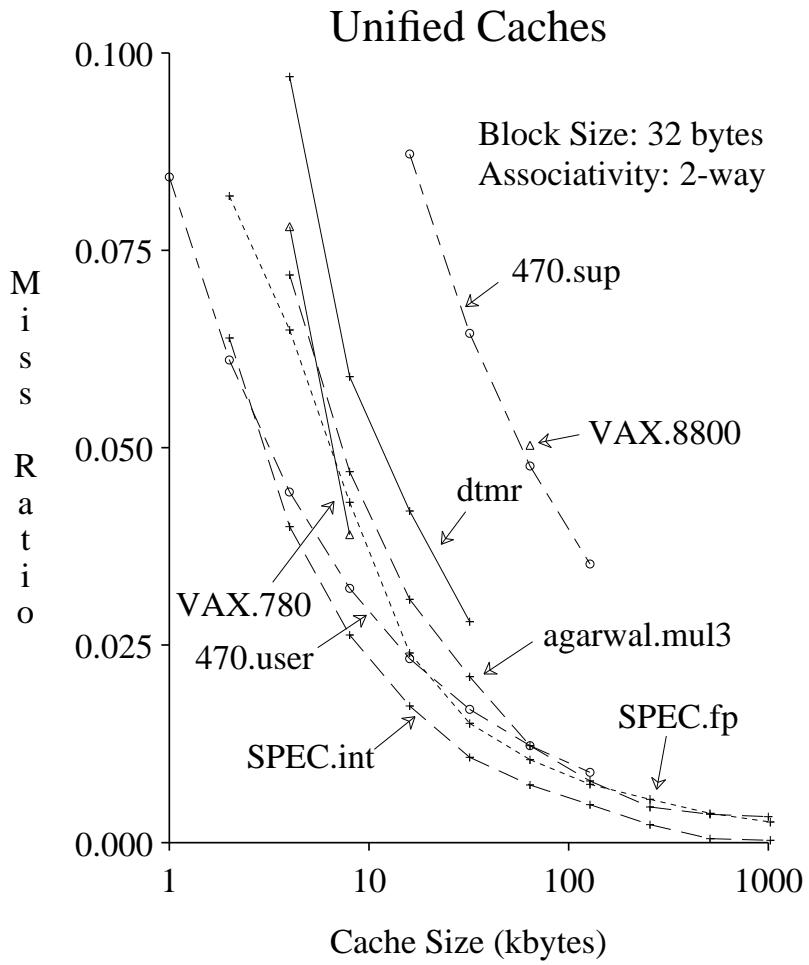


Figure 3: Unified cache miss ratios

4. Conclusions

The purpose of this study is two-fold: to show measurements of the cache performance of the SPEC benchmarks and to comment on the usefulness of those benchmarks for cache and memory system design. While the cache performance of the SPEC benchmarks varies from program to program, we have found that the floating-point benchmarks generally require much larger cache sizes relative to the integer benchmarks. The integer benchmarks use no more than 128 Kbytes of instruction and 128 Kbytes of data cache, while the floating-point programs can take advantage of data caches of a megabyte or more.

Comparisons with other studies show that the SPEC integer benchmarks have miss ratios much smaller than reported by any set of published measurements of hardware monitor results, those taken using a microcode tracer, or those from studies using very long traces. Miss ratios for the SPEC floating-point benchmarks seem consistent with previous measurements of user program miss ratios but are quite low relative to supervisor code miss ratios.

We note that there is no one unique workload or standard set of miss ratios; every environment will have its own workload and corresponding cache performance. From these measurements and comparisons, however, we conclude that miss ratios for the SPEC benchmarks could be considered representative of only a certain narrow environment - Unix workstations running user state CPU bound jobs as the single active user process. The integer benchmarks have very low miss ratios, and provide very little stress on the memory system. The floating point benchmarks provide reasonable measurements of memory system performance for user code, but are still much better behaved than commercial and timeshared workloads. The SPEC92 benchmarks are conspicuously lacking a significant operating system component, which affects their utility in two ways: miss ratios are very low, and the performance impacts of operating systems functions themselves are not tested.

We believe that an important aspect of the validity of a benchmark suite is that the benchmarks affect the memory system in a manner similar to that of the workload being represented. Our analysis above permits one to determine if the SPEC92 benchmark suite is suitable as a standard. Similar analysis should also be done for any subsequent release of the benchmarks.

5. Acknowledgments

Jeffrey Gee was with the Department of Electrical Engineering and Computer Sciences of University of California, Berkeley at the time this work was carried out. We would like to thank Gurindar Sohi and David Wood of University of Wisconsin-Madison for providing us with additional computational resources.

The material presented here is based on research supported in part by the National Science Foundation under grants MIP-8713274, MIPS-8957278, MIP-9116578, CCR-9117028 and CCR-9157366, by NASA under grant NCC2-550, by the State of California under the MICRO program, and by A.T.& T. Bell Laboratories, Cray Research Foundation, Digital Equipment Corporation, Intel Corporation, International Business Machines Corporation, Mitsubishi Electric Research Laboratories, Philips Laboratories/Signetics and Sun Microsystems.

References

- [Agar88] A. Agarwal, J. Hennessy, and M. Horowitz, "Cache Performance of Operating System and Multiprogramming Workloads," *ACM Trans. Comp. Sys.*, vol. 6, 4, November 1988, pp. 393-433.
- [Ande91] T.E. Anderson, H.M. Levy, B.N. Bershad, and E.D. Lazowska, "The Interaction of Architecture and Operating System Design," Proc. ASPLOS-IV, April, 1991, Santa Clara, CA, pp. 108-120.
- [Borg90] A. Borg, R.E. Kessler, and D.W. Wall, "Generation and Analysis of Very Long Address Traces," *Proc. 17th Int'l Symp. Comp. Arch.*, May, 1990, Seattle, WA, pp. 270-279.
- [Cmel91] R. M. Cmelik, S. I. Kong, D. R. Ditzel, and E. J. Kelly, "An Analysis of SPARC and MIPS Instruction Set Utilization on the SPEC benchmarks," Proc. ASPLOS-IV, April, 1991, Santa Clara, CA, pp. 290-302.
- [DEC91] "Pixie," DEC Ultrix manual page.
- [Clar83] D.W. Clark, "Cache Performance in the VAX-11/780," *ACM Trans. Comp. Sys.*, vol. 1, 1, February 1983, pp. 24-37.
- [Clar88] D.W. Clark, P.J. Bannon, J.B. Keller, "Measuring VAX 8800 Performance with a Histogram Hardware Monitor," *Proc. 15th Int'l Symp. Comp. Arch.*, May, 1988, Honolulu, HI, pp. 176-185.
- [Hill87] M.D. Hill, "Aspects of Cache Memory and Instruction Buffer Performance," Ph.D. Thesis, Univ. of California at Berkeley, Technical Report UCB/CSD 87/381, November 1987.
- [Hill89] Mark Hill and Alan Jay Smith, "Evaluating Associativity in CPU Caches," *IEEEETC*, 38, 12, December, 1989, pp. 1612-1630.
- [Hill] M.D. Hill, "Tycho," Unpublished UNIX-style manual page. The Tycho simulator is available from Prof. Mark Hill, Computer Sciences Dept., University of Wisconsin.
- [Hinn88] David Hinnant, "Accurate Unix Benchmarking: Art, Science or Black Magic?," *IEEE MICRO*, October, 1988, pp. 64-75.
- [Mogu91] J. C. Mogul, and Anita Borg, "The Effects of Context Switches on Cache Performance," Proc. ASPLOS-IV, April, 1991, Santa Clara, CA, pp. 75-84.
- [Pnev90] D.N. Pnevmatikatos, M.D. Hill, "Cache Performance of the Integer SPEC Benchmarks on a RISC," *Computer Architecture News*, vol. 18, 2, June 1990, pp. 53-68.
- [Pric89] Walter Price, "A Benchmark Tutorial," *IEEE MICRO*, October, 1989, pp. 28-43.
- [Saav90] Rafael H. Saavedra-Barrera, and Alan Jay Smith, "Performance Prediction by Benchmark and Machine Analysis," UC Berkeley Computer Science Division Technical Report UCB/CSD 90/607, December, 1990.
- [Saav92a] Rafael Saavedra-Barrera and Alan Jay Smith, "Analysis of Benchmark Characteristics and Benchmark Performance Prediction", Technical Report UCB/CSD-92-715, UC Berkeley Computer Science Division, December, 1992. Submitted for publication.
- [Saav92b] Rafael Saaavedra-Barrera, "CPU Performance Evaluation and Execution Time Prediction Using Narrow Spectrum Benchmarking", UC Berkeley Computer Science Technical Report UCB/CSD 92/684, February, 1992.
- [Sang84] John Sanguinetti, "Program Optimization for a Pipelined Machine: A Case Study," Proc. 1984 ACM Sigmetrics Conf. on Measurement and Modeling of Computer Systems, August, 1984, Cambridge, Mass., pp. 88-95.
- [Smit82] Alan Jay Smith, "Cache Memories," *Computing Surveys*, vol. 14, 3, September 1982.
- [Smit85] Alan Jay Smith, "Cache Evaluation and the Impact of Workload Choice," Proc. 12'th International Symposium on Computer Architecture, June 17-19, 1985, Boston, Mass, pp. 64-75.

- [Smit87] Alan Jay Smith, "Line (Block) Size Choice for CPU Cache Memories," *IEEE Trans. on Computers*, vol. C-36, 9, September 1987, pp. 1063-1075.
- [Spec91] SPEC newsletter, vol. 3, 4, 1991.