

- Using PAGE-1. Wiley-Interscience, New York, N.Y., 1972.
- Rein80a Rein, B. K. "A high-level approach to computer document formatting." In *Conf. Rec. 7th Annual ACM Symp. on Principles of Programming Languages* (Las Vegas, Nev., Jan. 1980), ACM, New York, 1980, pp. 24-31.
- Rein80b Rein, B. K., and Walker, J. H. *SCRIBE Introductory User's Manual*, 3rd ed., preliminary draft. Unilogic, Pittsburgh, 1980.
- Rein80c Rein, B. K. "Scribe: A Document Specification Language and Its Compiler." Ph.D. dissertation, Computer Science Dept., Carnegie-Mellon Univ., Pittsburgh, Pa., Oct. 1980. Also issued as Tech. Rep. CMU-CS-81-100.
- Rein81 Rein, B. K. "The Scribe document specification language and its compiler." In *Abstracts of the Presented Papers, Int. Conf. Research and Trends in Document Preparation Systems* (Lausanne, Switzerland, Feb. 1981), Swiss Institutes of Technology, Lausanne and Zurich, pp. 68-82.
- Rinc78 Rincenz, D. M. UNIX time sharing system: A retrospective. *Bell Syst. Tech. J.* 57, 6 (July-Aug. 1978), 1947-1969.
- Rob82a Rozengron, K. "ESP, a direct access editor: ESP user's guide." Tech. Note 134, Computer Science Lab., Univ. of Washington, Seattle, April 1981.
- Rob82b Rozengron, K. R. "ESP: A Direct Access Editor." Master's thesis, Univ. of Washington, Seattle, 1981.
- Sal76 Salzman, J. "Manuscript typing and editing: TYPSET, RUNOFF." In *The Compatible Time-Sharing System: A programmer's guide* 2nd ed., P. A. Christman (Ed.). The M.I.T. Press, Cambridge, Mass., 1966, sec. A1.9.01.
- Serv81 Servbold, J. "Xerox's 'Star'." *The Seybold Report* 10, 16 (April 27, 1981).
- Shaw80a Shaw, A. C. "A model for document preparation systems." Tech. Rep. 80-04-02, Dep. of Computer Science, Univ. of Washington, Seattle, April 1980.
- Shaw80b Shaw, A., Furuta, R., and Scofield, J. "Document formatting systems: Survey, concepts, and issues (Extended Abstract)." Tech. Rep. 80-10-02, Dep. of Comp. Sci., Univ. of Washington, Seattle, Oct. 1980. Also available in the *Abstracts of the Presented Papers, Int. Conf. Research and Trends in Document Preparation Systems* (Lausanne, Switzerland, Feb. 1981), Swiss Institutes of Technology, Lausanne and Zurich, pp. 47-52.
- Shreid80 Shreideman, B. *Software Psychology*. Wintrop, Cambridge, Mass., 1980.
- Shoch79 Shoch, J. F. "An overview of the programming language Smalltalk-72." *SIGPLAN Notices* (ACM) 14, 9 (Sept. 1979), 64-73.
- Smith75 Smith, D. C. "PYGMALION: A Creative Programming Environment." Ph.D. dissertation, Stanford Univ., Stanford, Calif., June 1975. Also issued as Stanford Artificial Intelligence Lab. Memo AIM-260 and as Computer Science Dep. Rep. STAN-CS-75-489.
- Smith82 Smith, D. C., Iken, C., Kimball, R., and Vesperlanc, B. *Designing the Star user interface. Byte* 7, 4 (April 1982), 242-252.
- Stryak80 Stryak, M. *The Joy of TEX: A personal guide to typesetting technical text by computer, Version 1.1*. American Mathematical Society, Providence, R.I., 1980.
- Stall80 Stallman, R. M. "EMACS manual for TENEX users." AI Memo 555, M.I.T. Artificial Intelligence Lab., Cambridge, Mass., Sept. 1980.
- Stall81 Stallman, R. M. "EMACS, the extensible, customizable self-documenting display editor." In *Proc. ACM SIGPLAN SIGOA Symp. on Text Manipulation, SIGPLAN Notices* (ACM) 16, 6 (June 1981), 147-156. Also available as *SIGOA Newsletter* (ACM) 2, 1&2 (Spring/Summer 1981), 147-156. This report is a revised version of AI Memo 519, M.I.T. Artificial Intelligence Lab., Cambridge, Mass., June 1979.
- Tesar81 Tesars, I. "PUB: The document compiler." Operating Note 70, Stanford Artificial Intelligence Project, Stanford, Calif., Sep. 1972.
- Thrac79 Thacker, C. P., McCarrioh, E. M., Lampron, B. W., Sproull, R. F., and Boggs, D. R. "Alto: A personal computer." Tech. Rep. CSL-79-11, Xerox Palo Alto Research Center, Palo Alto, Calif., Aug. 1979.
- Thom76 Thompson, K., and Rincenz, D. M. *UNIX Programmer's Manual*, 6th ed. Bell Telephone Laboratories, 1976, entry ROFF(1).
- Thom76 Van Dam, A., and Rice, D. E. On-line text editing: A survey. *ACM Comput. Surv.* 3, 3 (Sept. 1971), 83-114.
- VanL73 VanLehn, K. A. "SAL user manual." Rep. STAN-CS-73-373, Stanford Dep. of Comp. Sci., Stanford, Calif., July 1973. Also issued as Stanford Artificial Intelligence Lab. Memo AIM-204.
- VanW80 Van Wyk, C. J. "A Language for Typesetting Graphics." Ph.D. dissertation, Stanford Univ., Stanford, Calif., June 1980.
- VanW81 Van Wyk, C. J. A graphics typesetting language. In *Proc. ACM SIGPLAN SIGOA Symp. on Text Manipulation, SIGPLAN Notices* (ACM) 16, 6 (June 1981), 99-107. Also available as *SIGOA Newsletter* (ACM) 2, 1, 2 (Spring/Summer 1981), 99-107.
- Received November 1981; final revision accepted May 1982.

Cache Memories

ALAN JAY SMITH

University of California, Berkeley, California 94720

Cache memories are used in modern, medium and high-speed CPUs to hold temporarily those portions of the contents of main memory which are (believed to be) currently in use. Since instructions and data in cache memories can usually be referenced in 10 to 25 percent of the time required to access main memory, cache memories permit the execution rate of the machine to be substantially increased. In order to function effectively, cache memories must be carefully designed and implemented. In this paper, we explain the various aspects of cache memories and discuss in some detail the design features and trade-offs. A large number of original, trace-driven simulation results are presented. Consideration is given to practical implementation questions as well as to more abstract design issues.

Specific aspects of cache memories that are investigated include: the cache fetch algorithm (demand versus prefetch), the placement and replacement algorithms, line size, store-through versus copy-back updating of main memory, cold-start versus warm-start miss ratios, multichance consistency, the effect of input/output through the cache, the behavior of split data/instruction caches, and cache size. Our discussion includes other aspects of memory system architecture, including translation lookaside buffers.

Throughout the paper, we use as examples the implementation of the cache in the Amдах 470V/6 and 470V/7, the IBM 3081, 3033, and 370/168, and the DEC VAX 11/780. An extensive bibliography is provided.

Categories and Subject Descriptors: B.3.2 [Memory Structures]: Design Styles—cache memory; B.3.3 [Memory Structures]: Performance Analysis and Design Aids; C.O. [Computer Systems Organization]: General; C.4 [Computer Systems Organization]: Performance of Systems

General Terms: Design, Experimentation, Measurement, Performance

Additional Key Words and Phrases: Buffer memory, paging, prefetching, TLB, store-through, Amдах 470, IBM 3083, BIAS

INTRODUCTION

Definition and Rationale

Cache memories are small, high-speed buffer memories used in modern computer systems to hold temporarily those portions of the contents of main memory which are (believed to be) currently in use. Information located in cache memory may be accessed in much less time than that located in main memory (for reasons discussed throughout this paper). Thus, a central processing unit (CPU) with a cache memory needs to spend far less time waiting for

instructions and operands to be fetched and/or stored. For example, in typical large, high-speed computers (e.g., Amдах 470V/7, IBM 3033), main memory can be accessed in 300 to 600 nanoseconds; information can be obtained from a cache, on the other hand, in 50 to 100 nanoseconds. Since the performance of such machines is already limited in instruction execution rate by cache memory access time, the absence of any cache memory at all would produce a very substantial decrease in execution speed. Virtually all modern large computer sys-

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

CONTENTS

INTRODUCTION

Definition and Rationale
Overview of Cache Design
Cache Aspects

1. DATA AND MEASUREMENTS

1.1 Rationale
1.2 Trace-Driven Simulation
1.3 Simulation Evaluation
1.4 The Traces
1.5 Simulation Methods

2. ASPECTS OF CACHE DESIGN AND OPERATION

2.1 Cache Fetch Algorithm
2.2 Placement Algorithm
2.3 Line Size
2.4 Replacement Algorithm
2.5 Write-Through versus Copy-Back
2.6 Effect of Multiprogramming, Cold-Start, and Warm-Start
2.7 Multicache Consistency
2.8 Data/Instruction Cache
2.9 Virtual Address Cache
2.10 User/Supervisor Cache
2.11 Input/Output through the Cache
2.12 Cache Size
2.13 Cache Bandwidth, Data Path Width, and Access Latency
2.14 Multilevel Cache
2.15 Pipelining
2.16 Translation Lookaside Buffer
2.17 Translator
2.18 Memory-Biased Cache
2.19 Specialized Caches and Cache Components

3. DIRECTIONS FOR RESEARCH AND DEVELOPMENT

3.1 On-Chip Cache and Other Technology Advances
3.2 Multicache Consistency
3.3 Implementation Evaluation
3.4 Hit Ratio versus Size
3.5 TLB Design
3.6 Cache Parameters versus Architecture and Workload

APPENDIX. EXPLANATION OF TRACE NAMES AND ACKNOWLEDGMENTS

REFERENCES

tems have cache memories; for example, the Amdahl 470, the IBM 3081 [IBM82, Ren82, Gue82], 3033, 370/168, 360/195, the Univac 1100/80, and the Honeywell 66/80. Also, many medium and small size machines have cache memories; for example, the DEC VAX 11/780, 11/760 [ARMS81], and PDP-11/70 [Str86, Snow78], and the Apollo, which uses a Motorola 68000 microprocessor. We believe that within

the last four years, circuit speed and density will progress sufficiently to permit cache memories in one chip microcomputers. (On-chip addressable memory is planned for the Texas Instruments 99000 [Lar81, Erc81].) Even microcomputers could benefit substantially from an on-chip cache, since on-chip access times are much smaller than off-chip access times. Thus, the material presented in this paper should be relevant to almost the full range of computer architecture implementations.

The success of cache memories has been explained by reference to the "property of locality" [Denn72]. The property of locality has two aspects, temporal and spatial. Over short periods of time, a program distributes its memory references nonuniformly over its address space, and which portions of the address space are favored remain largely the same for long periods of time. This first property, called temporal locality, or locality by time, means that the information which will be in use in the near future is likely to be in use already. This property of behavior can be expected from programs in loops in which both data and instructions are reused. The second property, locality by space, means that portions of the address space which are in use generally consist of a fairly small number of individually contiguous segments of that address space. Locality by space, then, means that the local reference of the program in the near future are likely to be near the current local reference. This type of behavior can be expected from common knowledge of programs: related data items (variables, arrays) are usually stored together, and instructions are mostly executed sequentially. Since the cache's memory buffers segments of information that have been recently used, the property of locality implies that needed information is also likely to be found in the cache.

Optimizing the design of a cache memory generally has four aspects:

- (1) Maximizing the probability of finding a memory reference's target in the cache (the hit ratio),
- (2) minimizing the time to access information that is indeed in the cache (access time),
- (3) minimizing the delay due to a miss, and

- (4) minimizing the overheads of updating main memory, maintaining multicache consistency, etc.

(All of these have to be accomplished under suitable cost constraints, of course.) There is also a trade-off between hit ratio and access time. This trade-off has not been sufficiently stressed in the literature and it is one of our major concerns in this paper.

In this paper, each aspect of cache memories is discussed at length and, where available, measurement results are presented. In order for these detailed discussions to be meaningful, a familiarity with many of the aspects of cache design is required. In the remainder of this section, we explain the operation of a typical cache memory, and then we briefly discuss several aspects of cache memory design. These discussions are expanded upon in Section 2. At the end of this paper, there is an extensive bibliography in which we have attempted to cite all relevant literature. Not all of the items in the bibliography are referenced in the paper, although we have referred to items there as appropriate. The reader may wish in particular to refer to Bad79, Bars72, Gib87, and Kap73 for other surveys of some aspects of cache design. CLAR81 is particularly interesting as it discusses the design details of a real cache. (See also Lam80.)

Overview of Cache Design

Many CPUs can be partitioned, conceptually and sometimes physically, into three parts: the I-unit, the E-unit, and the S-unit. The I-unit (instruction) is responsible for instruction fetch and decode. It may have some local buffers for lookahead prefetching of instructions. The E-unit (execution) does most of what is commonly referred to as *executing* an instruction, and it contains the logic for arithmetic and logical operations. The S-unit (storage) provides the memory interface between the I-unit and E-unit. (IBM calls the S-unit the PSCF, or processor storage control function.)

The S-unit is the part of the CPU of primary interest in this paper. It contains several parts or functions, some of which are shown in Figure 1. The major component of the S-unit is the cache memory.

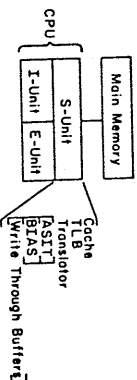


Figure 1. A typical CPU design and the S-unit.

There is usually a translator, which translates virtual to real memory addresses, and a TLB (translation lookaside buffer) which buffers (caches) recently generated (virtual address, real address) pairs. Depending on machine design, there can be an ASIT (address space identifier table), a BIAS (buffer invalidation address stack), and some write-through buffers. Each of these is discussed in later sections of this paper.

Figure 2 is a diagram of portions of a typical S-unit, showing only the more important parts and data paths, in particular the cache and the TLB. This design is typical of that used by IBM (in the 370/168 and 3033) and by Amdahl (in the 470 series). Figure 3 is a flowchart that corresponds to the operation of the design in Figure 2. A discussion of this flowchart follows.

The operation of the cache commences with the arrival of a virtual address, generally from the CPU, and the appropriate control signal. The virtual address is passed to both the TLB and the cache storage. The TLB is a small associative memory which maps virtual to real addresses. It is often organized as shown, as a number of groups (sets) of elements, each consisting of a virtual address and a real address. The TLB accepts the virtual page number, randomizes it, and uses that hashed number to select a set of elements. That set of elements is then searched associatively for a match to the virtual address. If a match is found, the corresponding real address is passed along to the comparator to determine whether the target line is in the cache. Finally, the replacement status of each entry in the TLB set is updated.

If the TLB does not contain the (virtual address, real address) pair needed for the translation, then the translator (not shown in Figure 2) is invoked. It uses the high-order bits of the virtual address as an entry into the segment and page tables for the

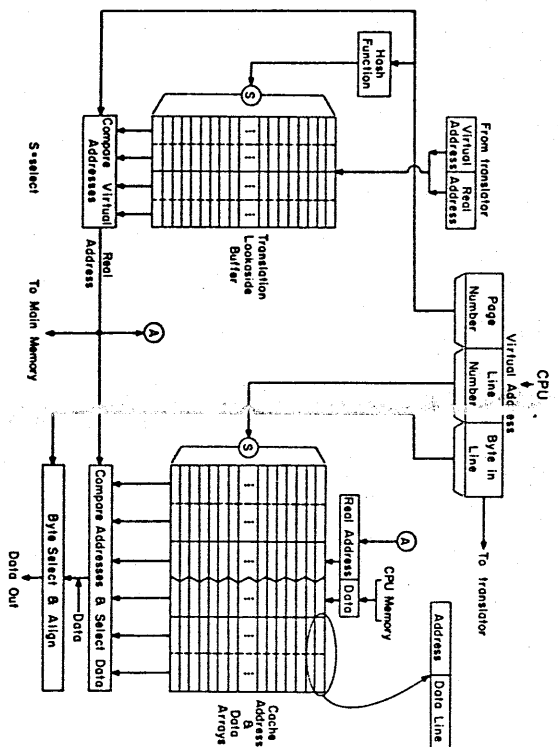


Figure 2. A typical cache and TLB design.

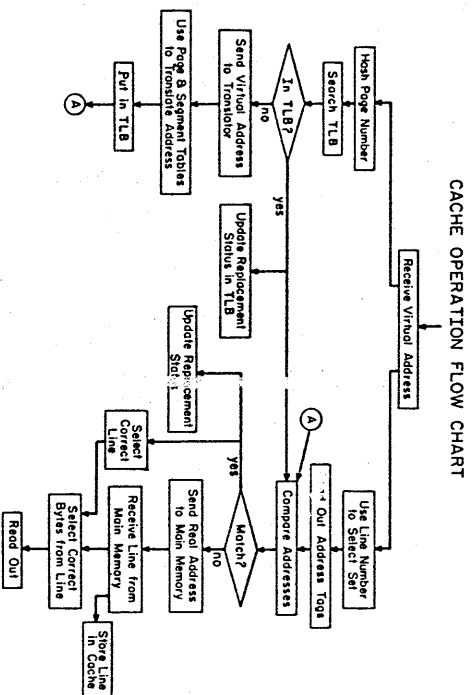


Figure 3. Cache operation flow chart.

process and then returns the address pair to the TLB (which retains it for possible future use), thus replacing an existing TLB entry.

The virtual address is also passed along initially to a mechanism which uses the middle part of the virtual address (the line number) as an index to select a set of entries in the cache. Each entry consists primarily of a real address tag and a line of data (see Figure 4). The line is the quantum of storage in the cache. The tags of the elements of all the selected set are read into a comparator and compared with the real address from the TLB. (Sometimes the cache stores the data and address tags together, as shown in Figures 2 and 4. Other times, the address tags and data are stored separately in the "address array" and "data array," respectively.) If a match is found, the line (or a part of it) containing the target locations is read into a shift register and the replacement status of the entries in the cache set are updated. The shift register is then shifted to select the target bytes, which are in turn transmitted to the source of the original data request.

If a miss occurs (i.e., address tags in the cache do not match), then the real address of the desired line is transmitted to the main memory. The replacement status information is used to determine which line to remove from the cache to make room for the target line. If the line to be removed from the cache has been modified, and main memory has not yet been updated with the modification, then the line is copied back to main memory; otherwise, it is simply deleted from the cache. After some number of machine cycles, the target line arrives from main memory and is loaded into the cache storage. The line is also passed to the shift register for the target bytes to be selected.

Cache Aspects

The cache description given above is both simplified and specific; it does not show design alternatives. Below, we point out some of the design alternatives for the cache memory.

Cache Fetch Algorithm. The cache fetch algorithm is used to decide when to bring information into the cache. Several possi-

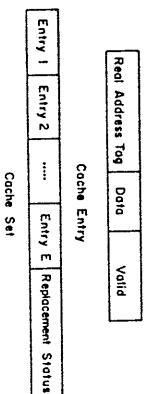


Figure 4. Structure of cache entry and cache set.

bilities exist: information can be fetched on demand (when it is needed) or prefetched (before it is needed). Prefetch algorithms attempt to guess what information will soon be needed and obtain it in advance. It is also possible for the cache fetch algorithm to omit fetching some information (selective fetch) and designate some information, such as shared writable code (semaphores), as unfetchable. Further, there may be no fetch-on-write in systems which use write-through (see below).

Cache Placement Algorithm. Information is generally retrieved from the cache associatively, and because large associative memories are usually very expensive and somewhat slow, the cache is generally organized as a group of smaller associative memories. Thus, only one of the associative memories has to be searched to determine whether the desired information is located in the cache. Each such (small) associative memory is called a set and the number of elements over which the associative search is conducted is called the set size. The placement algorithm is used to determine in which set a piece (line) of information will be placed. Later in this paper we consider the problem of selecting the number of sets, the set size, and the placement algorithm in such a set-associative memory.

Line Size. The fixed-size unit of information transfer between the cache and main memory is called the line. The line corresponds conceptually to the page, which is the unit of transfer between the main memory and secondary storage. Selecting the line size is an important part of the memory system design. (A line is also sometimes referred to as a block.)

Replacement Algorithm. When information is requested by the CPU from main memory and the cache is full, some information in the cache must be selected for

replacement. Various replacement algorithms are possible, such as FIFO (first in, first out), LRU (least recently used), and random. Later, we consider the first two of these.

Main Memory Update Algorithm. When the CPU performs a write (store) to memory, that operation can actually be reflected in the cache and main memories in a number of ways. For example, the cache memory can receive the write and the main memory can be updated when that line is replaced in the cache. This strategy is known as *copy-back*. Copy-back may also require that the line be fetched if it is absent from the cache (i.e., fetch-on-write). Another strategy, known as *write-through*, immediately updates main memory when a write occurs. Write-through may specify that if the information is in the cache, the cache be either updated or purged from main memory. If the information is not in the cache, it may or may not be fetched. The choice between copy-back and write-through strategies is also influenced by the need to maintain consistency among the cache memories in a tightly coupled multiprocessor system. This requirement is discussed later.

Cold-Start versus Warm-Start Miss Ratios and Multiprogramming. Most computer systems with cache memories are multiprogrammed: many processes run on the CPU, though only one can run at a time, and they alternate every few milliseconds. This means that a significant fraction of the cache miss ratio is due to loading data and instructions for a new process, rather than to a single process which has been running for some time. Miss ratios that are measured when starting with an empty cache are called *cold-start miss ratios*, and those that are measured from the time the cache becomes full are called *warm-start miss ratios*. Our simulation studies consider this multiprogramming environment.

User/Supervisor Cache. The frequent switching between user and supervisor state in most systems results in high miss ratios because the cache is often reloaded (i.e., cold-start). One way to address this is to incorporate two cache memories, and allow the supervisor to use one cache and

the other to use the other. Potentially, this could result in both the supervisor and the user programs more frequently finding upon initiation what they need in the cache.

Multicache Consistency. A multiprocessor system with multiple caches faces the problem of making sure that all copies of a given piece of information (which potentially could exist in every cache, as well as in the main memory) are the same. A modification of any one of these copies should somehow be reflected in all others. A number of solutions to this problem are possible. Three most popular solutions are essentially: (1) to transmit all stores to all caches and memories, so that all copies are updated; (2) to transmit the addresses of all stores to all other caches, and purge the corresponding lines from all other caches; or (3) to permit data that are writable (page or line flagged to permit modification) to be in only one cache at a time. A centralized or distributed directory may be used to control making and updating of copies.

Input/Output. Input/output (from and to I/O devices) is an additional source of references to information in memory. It is important that an output request stream reference the most current values for the information transferred. Similarly, it is also important that input data be immediately reflected in any and all copies of those lines in memory. Several solutions to this problem are possible. One is to direct the I/O stream through the cache itself (in a single processor system); another is to use a write-through policy and broadcast all writes so as to update or invalidate the target line wherever found. In the latter case, the channel accesses main memory rather than the cache.

Data/Instruction Cache. Another cache design strategy is to split the cache into two parts: one for data and one for instructions. This has the advantages that the bandwidth of the cache is increased and the access time (for reasons discussed later) can be decreased. Several problems occur: the overall miss ratio may increase, the two caches must be kept consistent, and self-modifying code and execute instructions must be accommodated.

Virtual versus Real Addressing. In computer systems with virtual memory, the cache may potentially be accessed either with a real address (real address cache) or with a virtual address (virtual address cache). If a virtual address (virtual address cache). If a real address is to be used, the virtual addresses generated by the processor must first be translated as in the example above (Figure 2); this is generally done by a TLB. The TLB is itself a cache memory which stores recently used address translation information, so that translation can occur quickly. Direct virtual address access is faster (since no translation is needed), but causes some problems. In a virtual address cache, inverse mapping (real to virtual address) is sometimes needed; this can be done by an RTB (reverse translation buffer).

Cache Size. It is obvious that the larger the cache, the higher the probability of finding the needed information in it. Cache sizes cannot be expanded without limit, however, for several reasons: cost (the most important reason in many machines, especially small ones), physical size (the cache must fit on the boards and in the cabinets), and access time. (The larger the cache, the slower it may become. Reasons for this are discussed in Section 2.12.) Later, we address the question of how large is large enough.

Multilevel Cache. As the cache grows in size, there comes a point where it may be usefully split into two levels: a small, high-level cache, which is faster, smaller, and more expensive per byte, and a larger, second-level cache. This two-level cache structure solves some of the problems that afflict caches when they become too large.

Cache Bandwidth. The cache bandwidth is the rate at which data can be read from and written to the cache. The bandwidth must be sufficient to support the proposed rate of instruction execution and I/O. Bandwidth can be improved by increasing the width of the data path, interleaving the cache and decreasing access time.

simulation. In this section, we explain the importance of this approach, and then discuss the presentation of our results.

One difficulty in providing definitive statements about aspects of cache operation is that the effectiveness of a cache memory depends on the workload of the computer system; further, to our knowledge, there has never been any (public) effort to characterize that workload with respect to its effect on the cache memory. Along the same lines, there is no generally accepted model for program behavior, and still less is there one for its effect on the uppermost level of the memory hierarchy. (But see AROR72 for some measurements, and LEM78 and LEM80, in which a model is used.)

For these reasons, we believe that it is possible for many aspects of cache design to make statements about relative performance only when those statements are based on trace-driven simulation or direct measurement. We have therefore tried throughout, when examining certain aspects of cache memories, to present a large number of simulation results and, if possible, to generalize from those measurements. We have also made an effort to locate and reference other measurement and trace-driven simulation results reported in the literature. The reader may wish, for example, to read WIND73, in which that author discusses the set of data used for his simulations.

1.2 Trace-Driven Simulation

Trace-driven simulation is an effective method for evaluating the behavior of a memory hierarchy. A trace is usually gathered by interpretively executing a program and recording every main memory location referenced by the program during its execution. (Each address may be tagged in any way desired, e.g., *instruction*, *fetch*, *data*, *fetch*, *data*, *store*.) One or more such traces are then used to drive a simulation model of a cache (or main) memory. By varying parameters of the simulation model, it is possible to simulate directly any cache size, placement, fetch or replacement algorithm, line size, and so forth. Programming techniques allow a range of values for many of these parameters to be measured simulta-

neously, during the same simulation run [GCS74, MAT70, SLUR72]. Trace-driven simulation has been a mainstay of memory hierarchy evaluation for the last 12 to 15 years; see BELA66 for an early example of this technique, or see PDM73. We assume only a single cache in the system, the one that we simulate. Note that our model does not include the additional buffers commonly found in the instruction decode and ALU portions of many CPUs.

In many cases, trace-driven simulation is preferred to actual measurement. Actual measurements require access to a computer and hardware measurement tools. Thus, if the results of the experiments are to be even approximately repeatable, standalone time is required. Also, if one is measuring an actual machine, one is unable to vary most (if any) hardware parameters. Trace-driven simulation has none of these difficulties; parameters can be varied at will and experiments can be repeated and reproduced precisely. The principal advantage of measurement over simulation is that it requires 1 to 0.1 percent as much running time and is thus very valuable in establishing a genuine, workload-based, actual level of performance (for validation). Actual workloads also include supervisor code, interrupts, context switches, and other aspects of workload behavior which are hard to imitate with traces. The results in this paper are mostly of the trace-driven variety.

1.3 Simulation Evaluation

There are two aspects to the performance of a cache memory. The first is access time. How long does it take to get information from or put information into the cache? It is very difficult to make exact statements about the effect of design changes on access time without specifying a circuit technology and a circuit diagram. One can, though, indicate trends, and we do that throughout this paper.

The second aspect of cache performance is the miss ratio. What fraction of all memory references attempt to access something which is not resident in the cache memory? Every such miss requires that the CPU wait until the desired information can be reached. Note that the miss ratio is a func-

tion not only of how the cache design affects the number of misses, but also of how the machine design affects the number of cache memory references. (A memory reference represents a cache access. A given instruction requires a varying number of memory references, depending on the specific implementation of the machine.) For example, a different number of memory references would be required if one word at a time were obtained from the cache than if two words were obtained at once. Almost all of our trace-driven studies assume a cache with a one-word data path (370 words = 4 bytes, PDP-11 word = 2 bytes). The WATEZ, WATFIV, FFT, and APL traces assume a two-word (eight-byte) data path. We measure the miss ratio and use it as the major figure of merit for most of our studies. We display many of these results as x/y plots of miss ratios versus cache size in order to show the dependence of various cache design parameters on the cache size.

1.4 The Traces

We have obtained 19 program address traces, 3 of them for the PDP-11 and the other 16 for the IBM 360/370 series of computers. Each trace is for a program developed for normal production use. (These traces are listed in the Appendix, with a brief description of each.) They have been used in groups to simulate multiprogramming; five such groups were formed. Two represent a scientific workload (WFFV, APL, WTX, FFT, and FGO1, FGO2, FGO3, FGO4), one a business (commercial) workload (CGO1, C702, CGO3, PGO2), one a miscellaneous workload, including compilations and a utility program (PGO1, CCOWP, FCOMP, IEBDG), and one a PDP-11 workload (ROFFAS, EDC, TRA2E). The miss ratio as a function of cache size is shown in Figure 5 for most of the traces; see SM79 for the miss ratios of the remaining traces. The miss ratios for each of the traces in Figure 5 are cold-start values based on simulations of 250,000 memory references for the IBM traces, and 333,333 for the PDP-11 traces.

1.5 Simulation Methods

Almost all of the simulations that were run used 3 or 4 traces and simulated multipro-

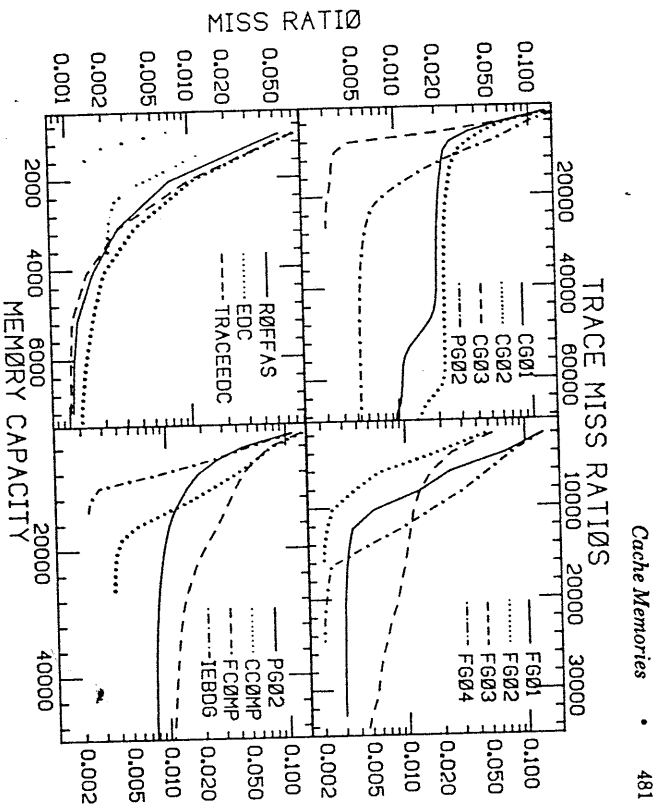


Figure 5. Individual trace miss ratios.

gramming by switching the trace in use every Q time-units (where Q was usually 10,000, a cache memory reference takes 1 time-unit, and a miss requires 10). Multiprogrammed simulations are used for two reasons: they are considered to be more representative of usual computer system operation than unprogrammed ones, and they also allow many more traces to be included without increasing the number of simulation runs. An acceptable alternative, though, would have been to use unprogrammed and purge the cache every Q memory references. A still better idea would have been to interleave user and supervisor code, but no supervisor traces were available.

All of the multiprogrammed simulations (i.e., Figures 6, 9-33) were run for one million memory references; thus approximately 250,000 memory references were used from each of the IBM 370 traces, and 333,333 from the PDP-11 traces.

The standard number of sets in the simulations was 64. The line size was generally 32 bytes for the IBM traces and 16 bytes for the PDP-11 traces.

2. ASPECTS OF CACHE DESIGN AND OPERATION

2.1 Cache Fetch Algorithm

2.1.1 Introduction

As we noted earlier, one of the two aims of cache design is to minimize the miss ratio. Part of the approach to this goal is to select a cache fetch algorithm that is very likely to fetch the right information, if possible, before it is needed. The standard cache fetch algorithm is demand fetching, by which a line is fetched when and if it is needed. Demand fetches cannot be avoided entirely, but they can be reduced if we can successfully predict which lines will be needed and fetch them in advance. A cache fetch algorithm which gets information be-

fore it is needed is called a prefetch algorithm.

Prefetch algorithms have been studied in detail in Smir78b. Below, we summarize those results and give one important extension. We also refer the reader to several other works [AICH76, BENN76, BERG78, ENGEL78, PERK80, and RAU76] for additional discussions of some of these issues.

We mention the importance of a technique known as *fetch bypass* or *load-through*. When a miss occurs, it can be rectified in two ways: either the line desired can be read into the cache, and the fetch then reinitiated (this was done in the original Amdahl 470V/6 [Smir78b]), or, better, the desired bytes can be passed directly from the main memory to the instruction unit, bypassing the cache. In this latter strategy, the cache is loaded, either simultaneously with the fetch bypass or after the bypass occurs. This method is used in the 470V/7, 470V/8, and the IBM 3033. (A wrap-around load is usually used [KROF80] in which the transfer begins with the bytes accessed and wraps around to the rest of the line.)

2.1.2 Prefetching

A prefetch algorithm must be carefully designed if the machine performance is to be improved rather than degraded. In order to show this more clearly, we must first define our terms. Let the *prefetch ratio* be the ratio of the number of lines transferred due to prefetches to the total number of program memory references. And let *transfer ratio* be the sum of the prefetch and miss ratios. There are two types of references to the cache: *actual* and *prefetch lookup*. Actual references are those generated by a source external to the cache, such as the rest of the CPU (I-unit, E-unit) or the channels. A prefetch lookup occurs when the cache interrogates itself to see if a given line is resident or if it must be prefetched. The ratio of the total accesses to the cache (actual plus prefetch lookup) to the number of actual references is called the *access ratio*.

There are costs associated with each of the above ratios. We can define these costs in terms of lost machine cycles per memory

reference. Let D be the penalty for a demand miss (a miss that occurs because the target is needed immediately) which arises from machine idle time while the fetch completes. The prefetch cost, P , results from the cache cycles used (and thus other lines unavailable) to bring in a prefetched line, used to move out (if necessary) a line replaced by a prefetch, and spent in delays while main memory modules are busy doing a pre-cache move-in and move-out. The access cost, A , is the penalty due to additional cache prefetch lookup accesses which interfere with the executing program's use of the cache. A prefetch algorithm is effective only if the following equation holds:

$$\begin{aligned} & \bullet \text{ miss ratio (demand)} \\ & > [D \bullet \text{ miss ratio (prefetch)} \\ & \quad + P \bullet \text{ prefetch ratio} \\ & \quad + A \bullet (\text{access ratio} - 1)] \quad (1) \end{aligned}$$

We should note also that the miss ratio when using prefetching may not be lower than the miss ratio for demand fetching. The problem here is cache *memory pollution*; prefetched lines may *pollute* memory by exelling other lines which are more likely to be referenced. This issue is discussed extensively and with some attempt at analysis in Smir78c; in Smir78b a number of experimental results are shown. We found earlier [Smir78b] that the major factor in determining whether prefetching is useful was the line size. Lines of 256 or fewer bytes (such as are commonly used in caches) generally resulted in useful prefetching; larger lines (or pages) made prefetching ineffective. The reason for this is that a prefetch to a large line brings in a great deal of information, much or all of which may not be needed, and removes an equally large amount of information, some of which may still be in use.

A prefetch algorithm has three major concerns: (1) when to initiate a prefetch, (2) which line(s) to prefetch, and (3) what replacement status to give the prefetched block. We believe that in cache memories, because of the need for fast hardware implementation, the only possible line to prefetch is the immediately sequential one; this type of prefetching is also known as *one block lookahead* (OBL). That is, if line

i is referenced, only line $i + 1$ is considered for prefetching. Other possibilities, which sometimes may result in a lower miss ratio, are not feasible for hardware implementation in a cache at cache speeds. Therefore, we consider only OBL.

If some lines in the cache have been referenced and others are resident only because they were prefetched, then the two types of lines may be treated differently with respect to replacement. Further, a prefetch lookup may or may not alter the replacement status of the line examined. In this paper we have made no distinction between the effect of a reference or a prefetch lookup on the replacement status of a line. That is, a line is moved to the top of the LRU stack for its set if it is referenced, prefetched, or is the target of a prefetch lookup; LRU is used for replacement for all prefetch experiments in this paper. (See Section 2.2.2) The replacement status of these three cases was varied in Smir78c, and in that paper it was found that such distinctions in replacement status had little effect on the miss ratio.

There are several possibilities for when to initiate a prefetch. For example, a prefetch can occur on instruction fetches, data reads and/or data writes, when a miss occurs, always, when the last rth of a line is accessed, when a sequential access pattern has already been observed, and so on. Prefetching when a sequential access pattern has been observed or when the last rth segment ($r = 1, 2, \dots$) of a line has been used is likely to be ineffective for reasons of timing: the prefetch will not be complete when the line is needed. In Smir78b we showed that limiting prefetches only to instruction accesses or only to data accesses is less effective than making all memory accesses eligible to start prefetches. See also BENN82.

It is possible to create prefetch algorithms or mechanisms which employ information not available within the cache memory. For example, a special instruction could be invented to initiate prefetches. No machine, to our knowledge, has such an instruction, nor have any evaluations been performed of this idea, and we are inclined to doubt its utility in most cases. A prefetch instruction that specified the transfer of

large amounts of information would run the substantial risk of polluting the cache with information that either would not be used for some time, or would not be used at all. If only a small amount of information were prefetched, the overhead of the prefetch might well exceed the value of the savings. However, some sophisticated versions of this idea might work. One such would be to make a record of the contents of the cache whenever the execution of a process was stopped, and after the process had been restarted, to restore the cache, or better, only its most recently used half. This idea is known as *working set restoration* and has been studied to some extent for paged main memories. The complexity of implementing it for cache makes it unlikely to be worthwhile, although further study is called for.

Another possibility would be to recognize when a base register is loaded by the processor and then to cause some number of lines (one, two, or three) following the loaded address to be prefetched [POWE80b, HOEY81a, HOEY81b]. Implementing this is easy, but architectural and software changes are required to ensure that the base registers are known or recognized, and modifications to them initiate prefetches. No evaluation of this idea is available, but a decreased miss ratio appears likely to result from its implementation. The effect could be very minor, though, and needs to be evaluated experimentally before any modification of current software or hardware is justified.

We consider three types of prefetching in this paper: (1) always prefetch, (2) prefetch on misses, and (3) tagged prefetch. *Always prefetch* means that on every memory reference, access to line i (for all i) implies a prefetch access for line $i + 1$. Thus the prefetch access ratio in this case is always 2.0. *Prefetch on misses* implies that a reference to a block i causes a prefetch to block $i + 1$ if and only if the reference to block i itself was a miss. Here, the access ratio is 1 + miss ratio. *Tagged prefetch* is a little more complicated, and was first proposed by GIND77. We associate with each line a single bit called the *tag*, which is set to one whenever the line is accessed by a program. It is initially zero and is reset to zero when

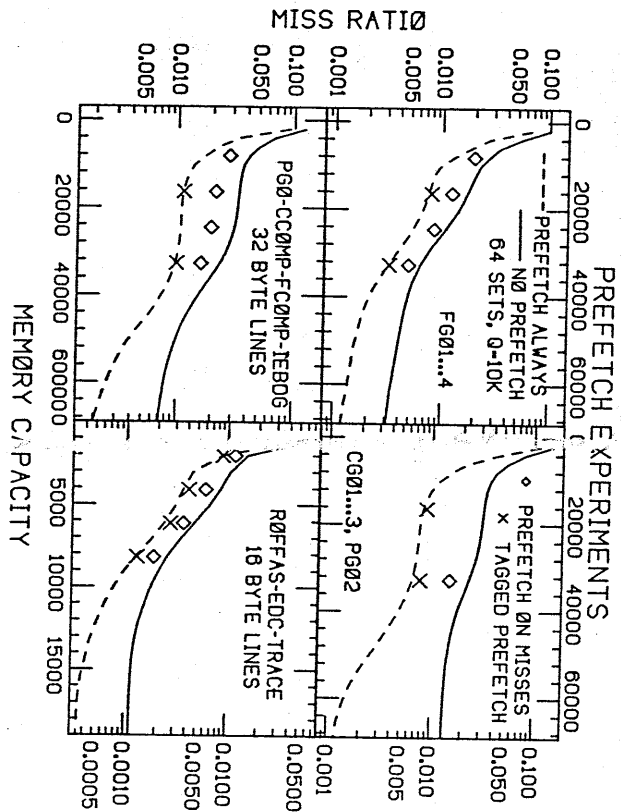


Figure 6. Comparison of miss ratios for two prefetch strategies and no prefetch.

the line is removed from the cache. Any line brought to the cache by a prefetch operation retains its tag of zero. When a tag changes from 0 to 1 (i.e., when the line is referenced for the first time after prefetching or is demand-fetched), a prefetch is initiated for the next sequential line. The idea is very similar to prefetching on misses only, except that a miss which did not occur because the line was prefetched (i.e., had there not been a prefetch, there would have been a miss to this line) also initiates a prefetch.

Two of these prefetch algorithms were tested in Smir78b: always prefetch and prefetch on misses. It was found that always prefetching reduced the miss ratio by as much as 75 to 80 percent for large cache memory sizes, while increasing the transfer ratio by 20 to 80 percent. Prefetching only on misses was much less effective; it produced only one half, or less, of the decrease in miss ratio produced by always prefetching. The transfer ratio, of course, also in-

creased by a much smaller amount, typically 10 to 20 percent.

The experiments in Smir78b, while very thorough, used only one set of traces and also did not test the tagged prefetch algorithm. To remedy this, we ran additional experiments; the results are presented in Figure 6. (In this figure, 32-byte lines are used in all cases except for 16-byte lines for the PDP-11 traces, the task switch interval $Q = 10K$, and there are 64 sets in all cases.) It can be seen that always prefetching cut the (demand) miss ratio by 50 to 90 percent for most cache sizes and tagged prefetch was almost equally effective. Prefetching on misses was less than half as good as always prefetching or tagged prefetch in reducing the miss ratio. These results are seen to be consistent across all five sets of traces used.

These experiments are confirmed by the results in Table 1. There we have tabulated the miss, transfer, and access ratios for the three prefetch algorithms considered, as

Table 1. Comparison of Three Prefetch Strategies*

Program traces	Memory size	Demand miss ratio	Always prefetch			Prefetch on misses			Tagged prefetch		
			Miss ratio	Access ratio	Transfer ratio	Miss ratio	Access ratio	Transfer ratio	Miss ratio	Access ratio	Transfer ratio
WFV, APL	16K	0.02162	0.00883	2.0	0.0297	—	—	—	0.00922	1.0233	0.01491
WTX, FFT1	32K	0.00910	0.00407	2.0	0.0152	0.00656	1.00656	0.01178	0.00405	1.0107	0.01275
ROFFAS, EDC	2K	0.0155	0.00867	2.0	0.0263	—	—	—	0.00873	1.0176	0.02245
TRACE	4K	0.00845	0.00356	2.0	0.0126	0.00593	1.0059	0.0107	0.00404	1.0098	0.01223
	6K	0.00470	0.00232	2.0	0.0085	—	—	—	0.00268	1.0059	0.00722
	8K	0.00255	0.00123	2.0	0.0047	0.00184	1.00184	0.00320	0.00127	1.0030	0.00362
CG01, CG02	16K	0.0341	0.00851	2.0	0.0391	0.01997	1.0197	0.03798	0.00893	1.0331	0.03892
CG03, PGO2	32K	0.0236	0.00670	2.0	0.0331	0.0153	1.0153	0.0288	0.00780	1.0274	0.03168
FG01, FG02	16K	0.01702	0.00736	2.0	0.0234	0.01234	1.01234	0.02239	0.00785	1.0194	0.0240
FG03, FG04	32K	0.00628	0.00314	2.0	0.0108	0.00489	1.00489	0.00668	0.00320	1.0077	0.00940
PG01, CCOMP	16K	0.0343	0.0112	2.0	0.0413	0.02087	1.02087	0.03928	0.01136	1.0333	0.0408
FCOMP, IEBD06	32K	0.0236	0.00960	2.0	0.0365	0.0163	1.0163	0.03010	0.0095	1.0286	0.0348

* Line size: 32 bytes for IBM 370 Traces, 16 byte lines for PDP-11 traces. Multiprogramming interval $Q = 10K$. Number of sets: 64.

well as the demand miss ratio for each of the sets of traces used and for a variety of memory sizes. We observe from this table that always prefetch and tagged prefetch are both very successful in reducing the miss ratio. Tagged prefetch has the significant additional benefit of requiring only a small increase in the access ratio over demand fetching. The transfer ratio is comparable for tagged prefetch and always prefetch.

It is important to note that taking advantage of the decrease in miss ratio obtained by these prefetch algorithms depends very strongly on the effectiveness of the implementation. For example, the Amdahl 470V/6-1 has a fairly sophisticated prefetch algorithm built in, but because of the design architecture of the machine, the benefit of prefetch cannot be realized. Although the prefetching cuts the miss ratio in this architecture, it uses too many cache cycles and interferes with normal program accesses to the cache. For that reason, prefetching is not used in the 470V/6-1 and is not available to customers. The more recent 470V/8, though, does contain a prefetch algorithm which is useful and improves machine performance. The V/8 cache prefetch algorithm prefetches (only) on misses, and was selected on the basis that it causes very little interference with normal machine operation. (Prefetch is implemented in the Dorado [Clair], but its success is not described.)

The prefetch implementation must at all times minimize its interference with regular machine functioning. For example, prefetch lookups should not block normal program memory accesses. This can be accomplished in three ways: (1) by instituting a second, parallel port to the cache, (2) by deferring prefetches until spare cache cycles are available, or (3) by not repeating recent prefetches. (Repeat prefetches can be eliminated by remembering the addresses of the last n prefetches in a small auxiliary cache. A potential prefetch could be tested against this buffer and not issued if found. This should cut the number of prefetch lookups by 80 to 90 percent for small n .)

The move-in (transfer of a line from main to cache memory) and move-out (transfer from cache to main memory) required by a

prefetch transfer can be buffered for performance during otherwise idle cache cycles. The main memory busy time engendered by a prefetch transfer seems unavoidable, but is not a serious problem. Also unavoidable is the fact that a prefetch may not be complete by the time the prefetched line is actually needed. This effect was examined in Smir78b and was found to be minor or although noticeable. Further comments and details of a suggested implementation are found in Smir78b.

We note briefly that it is possible to consider the successful use of prefetching as an indication that the line size is too small; prefetch functions much as a larger line size would. A comparison of the results in Figure 6 and Table 1 with those in Figure 15-21 shows that prefetching on misses for 12-byte lines gives slightly better results than doubling the line size to 64 bytes. Always prefetching and tagged prefetch are both significantly better than the larger line size without prefetching. Therefore, it would appear that prefetching has benefits in addition to those that it provides by simulating a larger line size.

2.2. Placement Algorithm

The cache itself is not a user-addressable memory, but serves only as a buffer for memory. Thus in order to locate an element in the cache, it is necessary to have some function which maps the main memory address into a cache location, or to search the cache associatively, or to perform some combination of these two. The placement algorithm determines the mapping function from main memory address to cache location.

The most commonly used form of placement algorithm is called set-associative mapping. It involves organizing the cache into S sets of E elements per set (see Figure 7). Given a memory address $r(i)$, a function $f(i)$ maps $r(i)$ into a set $s(i)$, so that $f(i) = s(i)$. The reason for this type of organization may be observed by letting either S or E become one. If S becomes one, then the cache becomes a fully associative memory. The problem is that the large number of lines in a cache would make a fully associative memory both slow and very expensive. (Our comments here apply

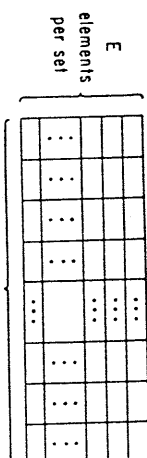


Figure 7. The cache is organized as S sets of E elements each.

to non-VLSI implementations. VLSI MOS facilitates broad associative searches.) Conversely, if E becomes one, in an organization known as direct mapping [CON69], there is only one element per set. (A more general classification has been proposed by HARD75.) Since the mapping function f is many to one, the potential for conflict in this latter case is quite high: two or more currently active lines may map into the same set. It is clear that, on the average, the conflict and miss ratio decline with increasing E , (as $S \cdot E$ remains constant), while the cost and access time increase. An effective compromise is to select E in the range of 2 to 16. Some typical values depending on certain cost and performance tradeoffs, are: 2 (Amdahl 470V/6-1, VAX 11/780, IBM 370/158-1), 4 (IBM 370/168-1, Amdahl 470V/8, Honeywell 66/80, IBM 370/158-3), 8 (IBM 370/168-3, Amdahl 470V/7), 16 (IBM 3033).

Another placement algorithm utilizes a sector buffer [CON68], as in the IBM 360/85. In this machine, the cache is divided into 16 sectors of 1024 bytes each. When a word is accessed for which the corresponding sector has not been allocated a place in the cache (the sector search being fully associative), a sector is made available (the LRU sector—see Section 2.4), and a 64-byte block containing the information referenced is transferred. When a word is referenced whose sector is in the cache but whose block is not, the block is simply fetched. The hit ratio for this algorithm is now generally known to be lower than that of the set-associative organization (Private Communication: F. Bookett) and hence we do not consider it further. (This type of design may prove appropriate for on-chip microprocessor caches, since the limiting factor in many microprocessor systems is

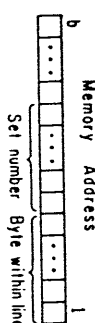


Figure 8. The set is selected by the low-order bits of the line number.

bus bandwidth. That topic is currently under study.)

There are two aspects to selecting a placement algorithm for the cache. First, the number of sets S must be chosen while $S \cdot E = M$ remains constant, where M is the number of lines in the cache. Second, the mapping function f , which translates a main memory address into a cache set, must be specified. The second question is most fully explored by Smir78a; we summarize those results and present some new experimental measurements below. A number of other papers consider one or both of these questions to some extent, and we refer to reader to those [CAM76, CON76, CON76, FUKU77, KAP73, LIP78, MAT71, STR76, THAK78] for additional information.

2.2.1 Set Selection Algorithm

Several possible algorithms are used or have been proposed for mapping an address into a set number. The simplest and most popular is known as bit selection, and is shown in Figure 8. The number of sets S is chosen to be a power of 2 (e.g., $S = 2^j$). If there are 2^j bytes per line, the j bits $1 \dots j$ select the byte within the line, and bits $j+1 \dots j+k$ select the set. Performing the mapping is thus very simple, since all that is required is the decoding of a binary quantity. Bit-selection is used in all com-

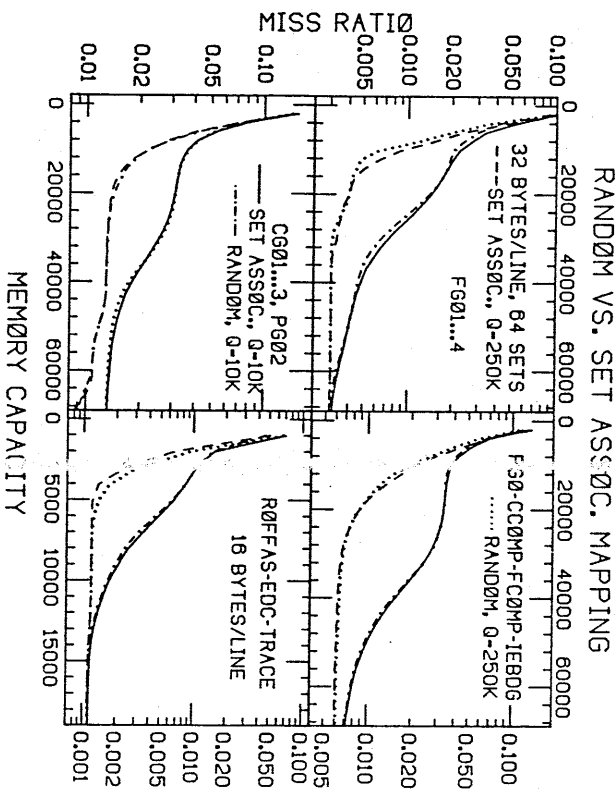


Figure 9. Comparison of miss ratios when using random or bit-selection set-associative mapping.

puters to our knowledge, including particularly all Amdahl and IBM computers.

Some people have suggested that because bit selection is not random, it might result in more conflict than a random algorithm, which employs a pseudorandom calculation (hashing) to map the line into a set. It is difficult to generate random numbers quickly in hardware, and the usual suggestion is some sort of *folding* of the address followed by exclusive or'ing of the bits. That is, if bits $j+1 \dots b$ are available for determining the line location, then these $b-j$ bits are grouped into k groups, and within each group an Exclusive Or is performed. The resulting k bits then designate a set. (This algorithm is used in TLBs—see Section 2.1.6.) In our simulations discussed later, we have used a randomizing function of the form $s(i) = a \cdot r(i) \bmod 2^k$.

Simulations were performed to compare random and set-associative mapping. The results are shown in Figure 9. (32 bytes is the line size used in all cases except the

PDP-11 traces, for which 16 bytes are used.) It can be observed that random mapping seems to have a small advantage in most cases, but that the advantage is not significant. Random mapping would probably be preferable to bit-selection mapping if it could be done equally quickly and inexpensively, but several extra levels of logic appear to be necessary. Therefore, bit selection seems to be the most desirable algorithm.

2.2.2 Set Size and the Number of Sets

There are a number of considerations in selecting values for the number of sets (S) and the set size (E). (We note that S and E are inversely related in the equation $S \cdot E = M$, where M is the number of lines in the cache ($M = 2^n$).) These considerations have to do with lookup time, expense, miss ratio, and addressing. We discuss each below.

The first consideration is that most cache memories (e.g., Amdahl, IBM) are addressed using the real address of the data,

although the CPU produces a virtual address. The most common mechanism for avoiding the time to translate the virtual address to a real one is to overlap the cache lookup and the translation operation (Figures 1 and 2). We observe the following: the only address bits that get translated in a virtual memory system are the ones that specify the page address, the bits that specify the byte within the page are invariant with respect to virtual memory translation. Let there be 2^p bytes per line and 2^k sets in the cache, as before. Also, let there be 2^p bytes per page. Then (assuming bit-selection mapping), $p-k$ bits are immediately available to choose the set. If $(p-k) \geq k$, then the set can be selected immediately, before translation, if $(p-k) < k$, then the search for the cache line can only be narrowed down to a small number $2^{k-(p-k)}$ of sets. It is quite advantageous if the set can be selected before translation (since the associative search can be started immediately upon completion of translation); thus there is a good reason to attempt to keep the number of sets less than or equal to $2^{(p-k)}$. (We note, though, that there is an alternative. The Amdahl 470V/6 has 256 sets, but only 6 bits immediately available for set selection. The machine reads out both elements of each of the four sets which could possibly be selected, then after the translation is complete, selects one of those sets before making the associative search. See also LEW80.)

Set size is just a different term for the scope of the associative search. The smaller the degree of associative search, the faster and less expensive the search (except, as noted above, for MOS VLSI). This is because there are fewer comparators and signal lines required and because the replacement algorithm can be simpler and faster (see Section 2.4). Our second consideration, expense, suggests that therefore the smaller the set size, the better. We repeat, though, that the set size and number of sets is inversely related. If the number of sets is less than or equal to $2^{(p-k)}$, then the set size is greater than or equal to $2^{(m-p-k)}$ lines.

The third consideration in selecting set size is the effect of the set size on the miss ratio. In [SMIT84a] we developed a model for this effect. We summarize the results of

that model here, and then we present some new experimental results.

A commonly used model for program behavior is what is known as the LRU Stack Model. (See COR73 for a more thorough explanation and some basic results.) In this model, the pages or lines of the program's address space are arranged in a list, with the most recently referenced line at the top of the list, and with the lines arranged in decreasing recency of reference from top to bottom. Thus, referencing a line moves it to the top of the list and moves all of the lines between the one referenced and the top line down one position. A reference to a line which is the i th line in the list (stack) is referred to as a stack distance of i . This model assumes that the stack distances are independently and identically drawn from a distribution $\{q(j)\}$, $j = 1, \dots, n$. This model has been shown not to hold in a formal statistical sense [LEW71, LEW73], but the author and others have used this model with good success in many modeling efforts.

Each set in the cache constitutes a separate associative memory. If each set is managed with LRU replacement, it is possible to determine the probability of referencing the k th most recently used item in a given set as a function of the overall LRU stack distance probability distribution $\{q(j)\}$. Let $p(i, S)$ be the probability of referencing the i th most recently referenced line in a set, given S sets. Then, we show that $p(i, S)$ may be calculated from the $\{q(i)\}$ with the following formula:

$$p(i, S) = \sum_{j=1}^n q(j) (1/S)^{j-1} (S - 1/S)^{j-1} \binom{j-1}{i-1}$$

Note that $p(i, 1) = q(i)$. In SMIT84a this model was shown to give accurate predictions of the effect of set size.

Experimental results are provided in Figures 10-14. In each of these cases, the number of sets has been varied. The rather curious shape of the curves (and the similarities between different plots) has to do with the task-switch interval and the fact that round-robin scheduling was used. Thus, when a program regained control of the processor, it might or might not, depending on the memory capacity, find any of its working set still in the cache.

VARY NUMBER OF SETS

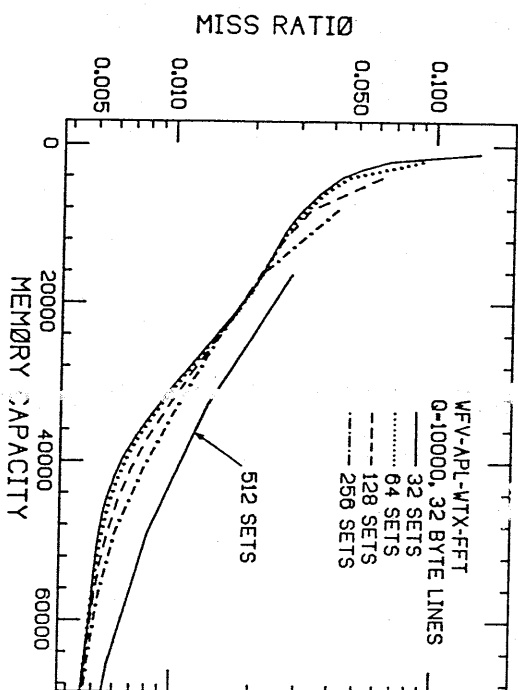


Figure 10. Miss ratios as a function of the number of sets and memory capacity.

VARY NUMBER OF SETS

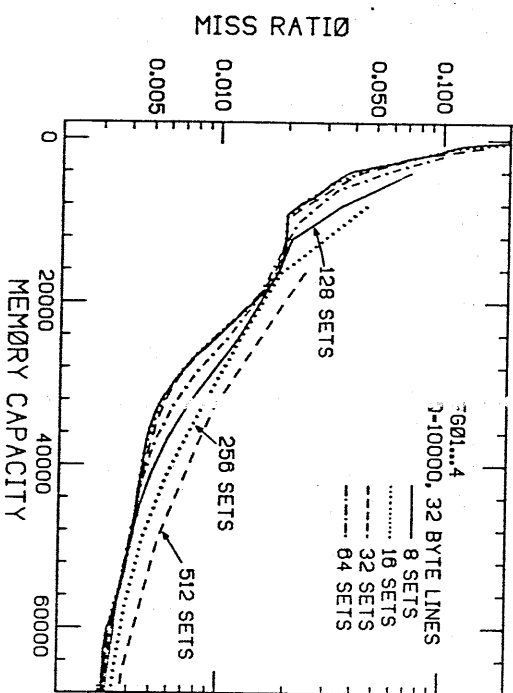


Figure 11. Miss ratios as a function of the number of sets and memory capacity.

VARY NUMBER OF SETS

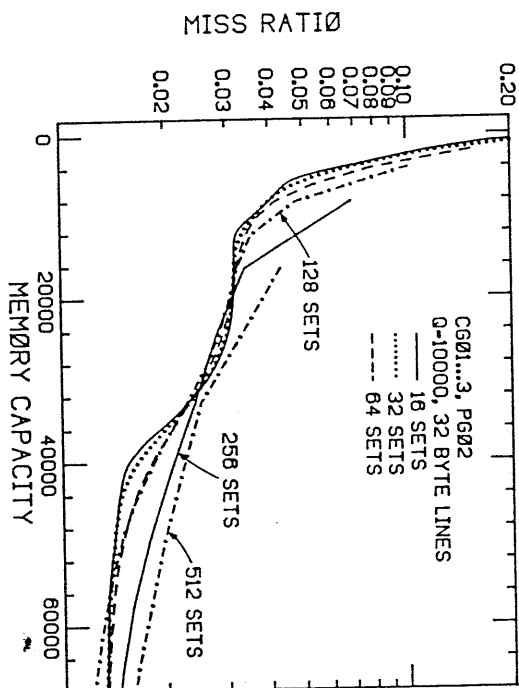


Figure 12. Miss ratios as a function of the number of sets and memory capacity.

VARY NUMBER OF SETS

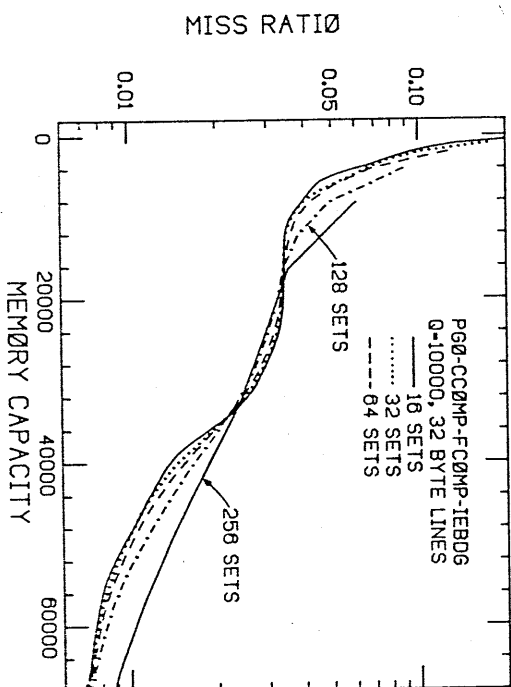


Figure 13. Miss ratios as a function of the number of sets and memory capacity.

VARY NUMBER OF SETS

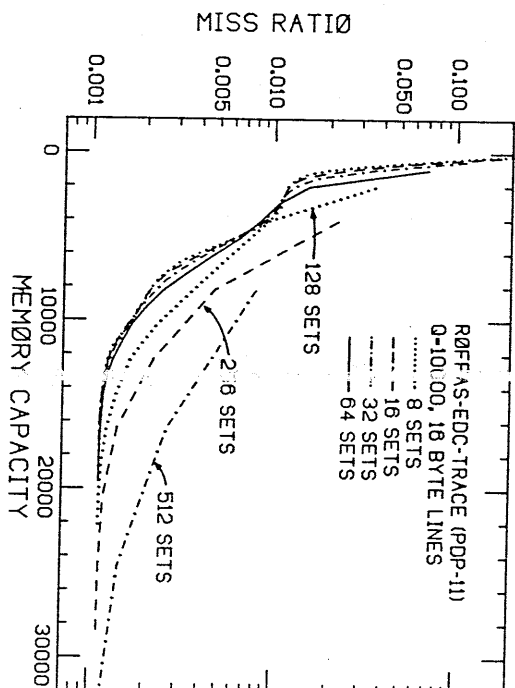


Figure 14. Miss ratios as a function of the number of sets and memory capacity.

Based on Figures 10-14, Figure 33, and the information given by Smitt78a, we believe that the minimum number of elements per set in order to get an acceptable miss ratio is 4 to 8. Beyond 8, the miss ratio is likely to decrease very little if at all. This has also been noted for TLB designs [SATT81]. The issue of maximum feasible set size also suggests that a set size of more than 4 or 8 will be inconvenient and expensive. The only machine known to the author with a set size larger than 8 is the IBM 3033 processor. The reason for such a large set size is that the 3033 has a 64-kbyte cache and 64-byte lines. The page size is 4096 bytes, which leaves only 6 bits for selecting the set, if the translation is to be overlapped. This implies that 16 lines are searched, which is quite a large number. The 3033 is also a very performance-oriented machine, and the extra expense is apparently not a factor.

Values for the set size (number of elements per set) and number of sets for a number of machines are as follows: Amdahl 470V/6 (2, 256), Amdahl 470V/7 (8, 128), Amdahl 470V/8 (4, 512), IBM 370/168-3 (8,

2.3 Line Size

One of the most visible parameters to be chosen for a cache memory is the line size. Just as with paged main memory, there are a number of trade-offs and no single criterion dominates. Below we discuss the advantages of both small and large line sizes. Additional information relating to this problem may be found in other papers [ALSA78, ANAC67, GINS67, KAPU73, METT71, METT70, and STRE76].

Small line sizes have a number of advantages. The transmission time for moving a set all line from main memory to cache is obviously shorter than that for a long line, at least if the machine has to wait for the full transmission time, short lines are better. (A high-performance machine will use fetch by pass; see Section 2.1.1.) The small line is less likely to contain unneeded information; only a few extra bytes are brought in along

VARY LINE SIZE

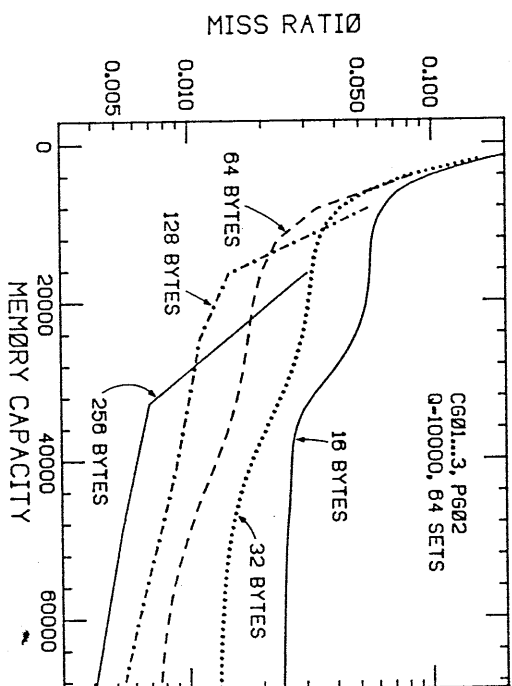


Figure 15. Miss ratios as a function of the line size and memory capacity.

with the actually requested information. The data width of main memory should usually be at least as wide as the line size, since it is desirable to transmit an entire line in one main memory cycle time. Main memory width can be expensive, and short lines minimize this problem.

Large line sizes, too, have a number of advantages. If more information in a line is actually being used, fetching it all at one time (as with a long line) is more efficient. The number of lines in the cache is smaller, so there are fewer logic gates and fewer storage bits (e.g., LRU bits) required to keep and manage address tags and replacement status. A larger line size permits fewer elements/set in the cache (see Section 2.2), which minimizes the associative search logic. Long lines also minimize the frequency of "line crossers," which are requests that span the contents of two lines. Thus in most machines, this means that two separate fetches are required within the cache (this is invisible to the rest of the machine.)

Note that the advantages cited above for

both long and short lines become disadvantages for the other.

Another important criterion for selecting a line size is the effect of the line size on the miss ratio. The miss ratio, however, only tells part of the story. It is inevitable that longer lines make processing a miss somewhat slower (no matter how efficient the overlapping and buffering), so that translating a miss ratio into a measure of machine speed is tricky and depends on the details of the implementation. The reader should bear this in mind when examining our experimental results.

Figures 15-21 show the miss ratio as a function of line size and cache size for five different sets of traces. Observe that we have also varied the multiprogramming quantum time Q . We do so because the miss ratio is affected by the task-switch interval nonuniformly with line size. This nonuniformity occurs because long line sizes load up the cache more quickly. Consider two cases. First, assume that most cache misses result from task switching. In this case, long lines load up the cache more

VARY LINE SIZE

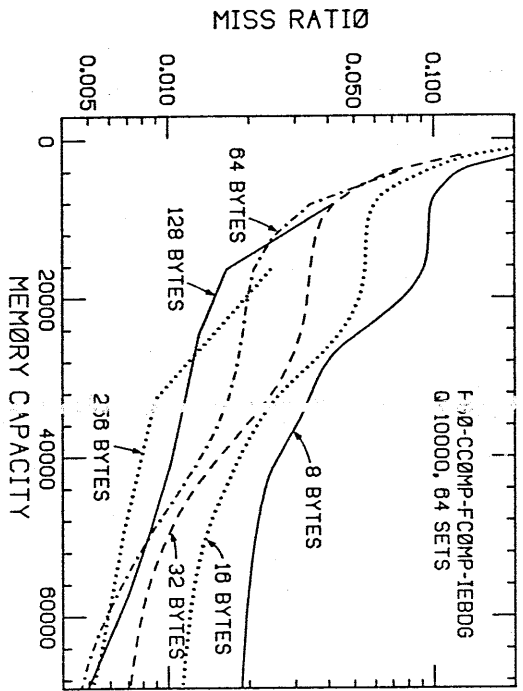


Figure 16. Miss ratios as a function of the line size and memory capacity.

VARY LINE SIZE

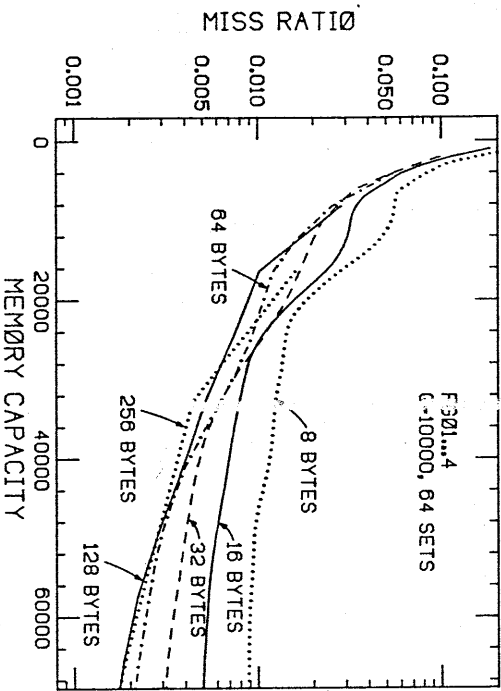


Figure 17. Miss ratios as a function of the line size and memory capacity.

VARY LINE SIZE

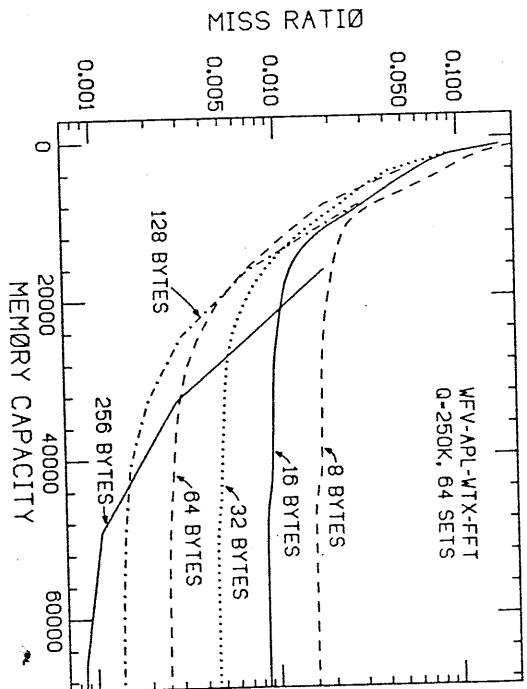


Figure 18. Miss ratios as a function of the line size and memory capacity.

VARY LINE SIZE

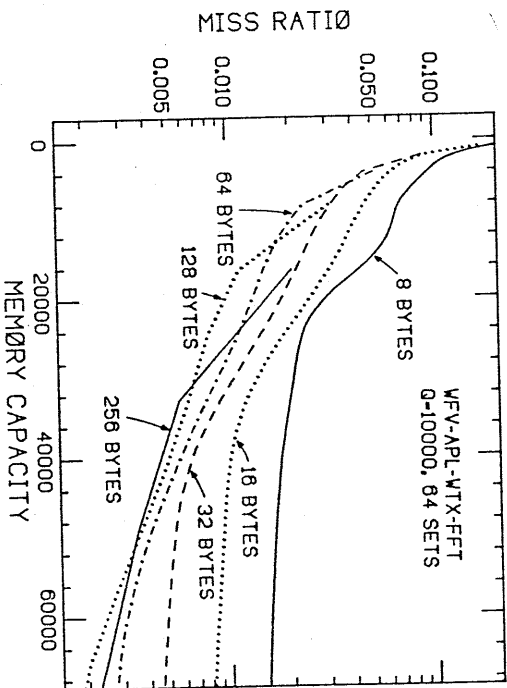


Figure 19. Miss ratios as a function of the line size and memory capacity.

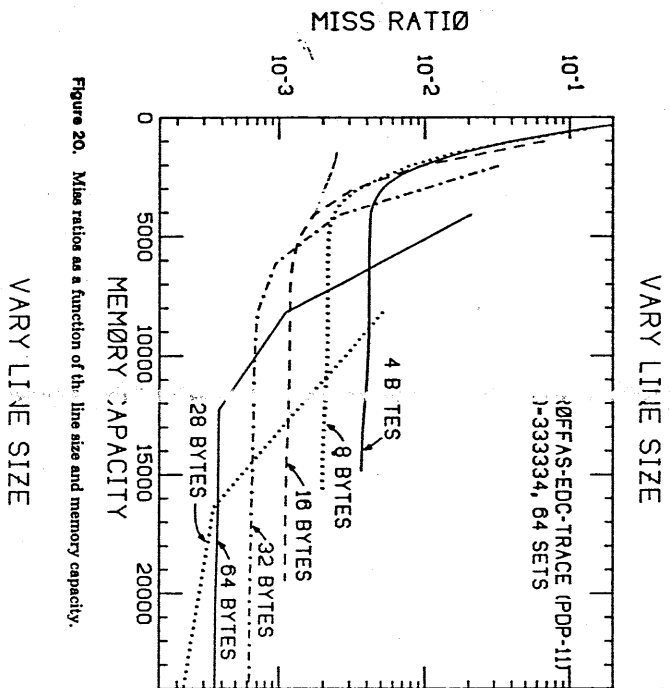


Figure 20. Miss ratios as a function of the line size and memory capacity.

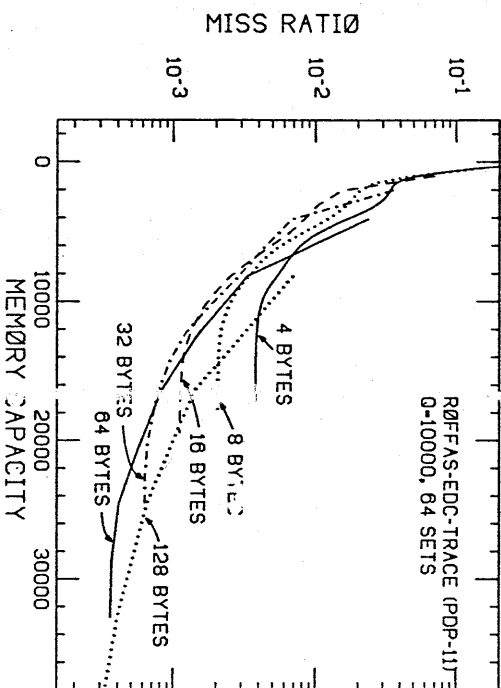


Figure 21. Miss ratios as a function of the line size and memory capacity.

quickly than small ones. Conversely, assume that most misses occur in the steady state; that is, that the cache is entirely full most of the time with the current process and most of the misses occur in this state. In this latter case, small lines cause less memory pollution and possibly a lower miss ratio. Some such effect is evident when comparing Figures 18 and 19 with Figures 20 and 21, but explanation is required. A quantum of 10,000 results not in a program finding an empty cache, but in it finding some residue of its previous period of activity (since the degree of multiprogramming is only 3 or 4); thus small lines are relatively more advantageous in this case than one would expect.

The most interesting number to be gleaned from Figures 15-21 is the line size which causes the minimum miss ratio for each memory size. This information has been collected in Table 2. The consistency displayed there for the 360/370 traces is surprising; we observe that one can divide the cache size by 128 or 256 to get the minimum miss ratio line size. This rule does not apply to the PDP-11 traces. Programs written for the PDP-11 not only use a different instruction set, but they have been written to run in a small (64K) address space. Without more data, generalizations from Figures 20 and 21 cannot be made.

In comparing the minimum miss ratio line sizes suggested by Table 2 and the offerings of the various manufacturers, one notes a discrepancy. For example, the IBM 168-1 (32-byte line, 16K buffer) and the 3033 (64-byte line, 64K buffer) both have surprisingly small line sizes. The reason for this is almost certainly that the transmission time for longer lines would induce a performance penalty, and the main memory data path width required would be too large and therefore too expensive.

Kumar [Kuma79] also finds that the line sizes in the IBM 3033 and Amdahl 470 are too small. He creates a model for the working set size w of a program, of the form $w(k) = c/k^a$, where k is the block size, and c and a are constants. By making some convenient assumptions, Kumar derives from this an expression for the miss ratio as a function of the block size. Both expressions are verified for three traces, and a is measured to be in the range of 0.45 to 0.85 over the

Table 2. Line Size (in bytes) Giving Minimum Miss Ratio, for Given Memory Capacity and Traces

Traces	Quantum	Memory size (bytes)	Minimum miss ratio line size (bytes)
CGO1	10,000	4	32
CGO2		8	64
CGO3		16	128
PGO2		32	256
PGO	10,000	64	256
CCOMP		4	32
FCOMP		8	64
IEBDG		16	128
IEBDG		32	256
PGO1	10,000	64	256
PGO2		4	32
PGO3		8	64
PGO4		16	128
WFO		32	256
WFO	10,000	64	256
APL		4	32
WTX		8	64
FTT		16	128
FTT		32	256
FTT	250,000	64	256
FTT		4	32
FTT		8	64
FTT		16	128
FTT		32	256
ROFFAS	10,000	2	16
EDC		4	32
TRACE		8	64
TRACE	333,333	16	128
TRACE		2	16
TRACE		4	32
TRACE		8	64
TRACE		16	128
TRACE		32	256
TRACE		64	256

three traces and various working set window sizes. He then found that for those machines, the optimum block size lies in the range 64 to 256 bytes.

It is worth considering the relationship between prefetching and line size. Prefetching can function much as a larger line size would. In terms of miss ratio, it is usually even better; although a prefetched line that is not being used can be swapped out, a half of a line that is not being used cannot be removed independently. Comparisons between the results in Section 2.1 and this section show that the performance improvement from prefetching is significantly larger than that obtained by doubling the line size.

Line sizes in use include: 128 bytes (IBM 3081 [IBM82]), 64 bytes (IBM 3033), 32 bytes (Amdahl 470s, Intel AS/6, IBM 370/

168), 16 bytes (Honeywell 66/60 and 66/80), 8 bytes (DEC VAX 11/780), 4 bytes (PDP-11/70).

2.4 Replacement Algorithm

2.4.1 Classification

In the steady state, the cache is full, and a cache miss implies not only a fetch but also a replacement; a line must be removed from the cache. The problem of replacement has been studied extensively for paged main memories (see Smir78d for a bibliography), but the constraints on a replacement algorithm are much more stringent for a cache memory. Principally, the cache replacement algorithm must be implemented entirely in hardware and must execute very quickly, so as to have no negative effect on processor speed. The set of feasible solutions is still large, but many of them can be rejected on inspection.

The usual classification of replacement algorithms groups them into usage-based versus non-usage-based, and fixed-space versus variable-space. Usage-based algorithms take the record of use of the line (or page) into account when doing replacement; examples of this type of algorithm are LRU (least recently used) [COFF73] and Working Set [DENN68]. Conversely, non-usage-based algorithms make the replacement decision on some basis other than and not related to usage; FIFO and Rand (random or pseudorandom) are in this class. (FIFO could arguably be considered usage-based, but since reuse of a line does not improve the replacement status of that line, we do not consider it as being such.) Fixed-space algorithms assume that the amount of memory to be allocated is fixed; replacement simply consists of selecting a specific line. If the algorithm varies the amount of space allocated to a specific process, it is known as a variable-space algorithm, in which case, a fetch does not imply a replacement, and a swap-out can take place without a corresponding fetch. Working Set and Page Fault Frequency [CHU76] are variable-space algorithms.

The cache memory is fixed in size, and it is usually too small to hold the working set of more than one process (although the 470V/8 and 3033 may be exceptions). For

this reason, we believe that variable-space algorithms are not suitable for a cache memory. To our knowledge, no variable-space algorithm has ever been used in a cache memory.

It should also be clear that in a set-associative memory, replacement must take place in the same set as the fetch. A line is being added to a given set because of the fetch, and thus a line must be removed. Since a line maps uniquely into a set, the replaced line in that set must be entirely removed from the cache.

The set of acceptable replacement algorithms is thus limited to fixed-space algorithms as executed within each set. The basic candidates are LRU, FIFO, and Rand. It is our experience (based on prior experiments and on material in the literature) that non-usage-based algorithms all yield comparable performance. We have chosen FIFO with a set as our example of a non-usage-based algorithm.

2.4.2 Comparisons

Comparisons between FIFO and LRU appear in Table 3, where we show results based on each set of traces for varying memory sizes, quantum sizes, and set number. We found (averaging over all of the numbers there) that FIFO yields a miss ratio approximately 12 percent higher than LRU, although the ratios of FIFO to LRU miss ratio range from 0.96 to 1.38. This 12 percent difference is significant in terms of performance, and LRU is clearly a better choice if the cost of implementing LRU is small and the implementation does not slow down the machine. We note that in minicomputers (e.g., PDP-11) cost is by far the major criterion; consequently, in such systems, this performance difference may not be worthwhile. The interested reader will find additional performance data and discussion in other papers [CHIA75, FURN78, GIBS67, LEE69, and STRE76].

2.4.3 Implementation

It is important to be able to implement LRU cheaply, and so that it executes quickly; the standard implementation in software using linked lists is unlikely to be either cheap or fast. For a set size of two,

Table 3. Miss Ratio Comparison for FIFO and LRU (within set) Replacement

Memory size (bytes)	Quantum size	Number of sets	Miss ratio LRU	Miss ratio FIFO	Ratio FIFO/ LRU	Traces
16	10K	64	0.02162	0.02254	1.04	WVU
32	10K	64	0.00910	0.01143	1.26	APL
16	250K	64	0.00868	0.01036	1.19	WTX
32	250K	64	0.00523	0.00548	1.05	FFT
16	10K	256	0.02171	0.02235	1.03	
32	10K	256	0.01067	0.01067	1.12	ROFPAS
4	10K	64	0.00845	0.00947	1.35	EBC
8	10K	64	0.00173	0.00344	1.94	TRACE
4	333K	64	0.00173	0.00214	1.00	
8	333K	64	0.00120	0.00120	1.00	
4	10K	256	0.02175	0.02175	0.998	
8	10K	256	0.0218	0.02175	1.08	
4	10K	256	0.00477	0.00514	1.00	
8	10K	256	0.01624	0.01624	1.00	
4	333K	256	0.00155	0.00159	1.03	CGO1
8	333K	256	0.00135	0.00139	1.38	CGO2
128	10K	64	0.01147	0.01103	0.96	CGO3
64	10K	64	0.01461	0.01694	1.30	PGO2
64	10K	256	0.01098	0.01171	1.08	PGO1
128	10K	256	0.01702	0.01872	1.10	FGO2
16	10K	64	0.00628	0.00687	1.38	FGO3
32	10K	64	0.01888	0.01934	1.02	FGO4
16	10K	256	0.00857	0.00946	1.10	PGO1
32	10K	256	0.03428	0.03496	1.02	PGO1
16	10K	64	0.02356	0.02543	1.08	CCOMP
32	10K	64	0.03540	0.03644	1.08	FCOMP
16	10K	256	0.03540	0.03644	1.06	IEBDG
32	10K	256	0.02394	0.02534	1.116	

only a hot/cold (toggle) bit is required. More generally, replacement in a set of E elements can be effectively implemented with $E(E-1)/2$ bits of status. (We note that $\log_2 E!$ bits of status are the theoretical minimum.) One creates an upper-left triangular matrix (without the diagonal, that is, $i < j$) which we will call R and refer to as $R(i, j)$. When line i is referenced, row i of $R(i, j)$ is set to 1, and column j of $R(i, j)$ is set to 0. The LRU line is the one for which the row is entirely equal to 0 (for those bits in the row, the row may be empty) and for which the column is entirely 1 (for all the bits in the column; the column may be empty). This algorithm can be easily implemented in hardware, and executes rapidly. See MARU76 for an extension and MARU76 for an alternative.

The above algorithm requires a number of LRU status bits that increases with the square of the set size. This number is acceptable for a set size of 4 (470V/8, IteL AS/6), marginal for a set size of eight (470V/7), and unacceptable for a set size of 16. For that reason, IBM has chosen to implement

approximations to LRU in the 370/168 and the 3033. In the 370/168-3 [IBM75], the set size is 8, with the 8 lines grouped in 4 pairs. The LRU pair of lines is selected, and then the LRU block of the pair is the one used for replacement. This algorithm requires only 10 bits, rather than the 28 needed by the full LRU. A set size of 16 is found in the 3033 [IBM78]. The 16 lines that make up a set are grouped into four groups of two pairs of two lines. The line to be replaced is selected as follows: (1) find the LRU group of four lines (requiring 6 bits of status), (2) find the LRU pair of the two pairs (1 bit per group, thus 4 more bits), and (3) find the LRU line of that pair (1 bit per pair, thus 8 more bits). In all, 18 bits are used for this modified LRU algorithm, as opposed to the 120 bits required for a full LRU. No experiments have been published comparing these modified LRU algorithms with genuine LRU, but we would expect to find no measurable difference.

Implementing either FIFO or Rand is much easier than implementing LRU. FIFO is implemented by keeping a modulo

Table 4. Percentage of Memory Reference¹ That Are Reads, Writes, and Instruction Fetches for Each Trace²

Trace	Partial trace			Full trace		
	Data read	Data write	IFETCH	Data read	Data write	IFETCH
WATFIV	23.3	16.54	60.12	—	16.89	—
APL	21.5	8.2	70.3	—	8.90	—
WATLEX	24.5	9.07	66.4	—	7.84	—
FRTI	23.4	7.7	68.9	—	7.59	—
ROFAS	37.5	4.96	57.6	—	6.4	—
EDC	30.4	10.3	59.2	29.8	11.0	59.2
TRACE	47.9	10.2	41.9	48.6	10.0	41.3
CGOI	41.5	34.2	24.3	42.07	34.19	23.74
CGO2	41.1	32.4	26.5	36.92	15.42	47.66
CGO3	37.7	22.5	39.8	37.86	22.55	39.59
PGO2	31.6	13.4	55.1	30.36	12.77	56.87
PGO1	29.9	17.6	52.6	30.57	11.26	58.17
PGO2	30.6	5.72	63.7	32.54	10.16	57.30
PGO3	30.0	12.8	57.2	30.60	13.25	56.15
PGO4	28.5	17.2	54.3	28.38	17.29	54.33
PGO1	29.7	19.8	50.5	28.68	16.83	54.39
CCOMP1	30.8	9.91	59.3	33.42	17.10	49.47
PCOMP1	29.5	20.7	50.0	30.80	16.51	53.68
IEBDG	39.3	28.1	32.7	39.3	28.2	32.5
Average	32.0	16.96	52.02	34.65	14.80	49.38
Stand. Dev.	7.1	8.7	13.4	6.90	7.21	10.56

¹ Partial trace results are for first 250,000 memory references for IBM 370 traces, and 333,333 memory references for PDP-11 traces. Full trace results refer to entire length of memory address trace (one to ten million memory percent reference).

E (*E* elements/set) counter for each set; it is incremented with each replacement and points to the next line for replacement. Rand is simpler still. One possibility is to use a single modulo *E* counter, incremented in a variety of ways: by each clock cycle, each memory reference, or each replacement anywhere in the cache. Whenever a replacement is to occur, the value of the counter is used to indicate the replaceable line within the set.

2.5 Write-Through versus Copy-Back

When the CPU executes instructions that modify the contents of the current address space, those changes must eventually be reflected in main memory; the cache is only a temporary buffer. There are two general approaches to updating main memory: stores can be immediately transmitted to main memory (called write-through or store-through), or stores can initially only modify the cache, and can later be reflected in main memory (copy-back). There are issues of performance, reliability, and complexity in making this choice; these issues are discussed in this section. Further information can be found in the literature

[Agia77a, Bell74, Pohl75, and Rist77]. A detailed analysis of some aspects of this problem is provided in Smrt79 and Yen81.

To provide an empirical basis for our discussion in this section, we refer the reader to Table 4. There we show the percentage of memory references that resulted from data reads, data writes, and instruction fetches for each of the traces used in this paper. The leftmost three columns show the results for those portions of the traces used throughout this paper; that is, the 70 traces were run for the first 250,000 memory references and the PDP-11 traces for 33,333 memory references. When available, the results for the entire trace (1 to 10 million memory references) are shown in the rightmost columns. The overall average shows 16 percent of the references were writes, but the variation is wide (5 to 34 percent) and the values observed are clearly very dependent on the source language and on the machine architecture. In Smrt79 we observed that the fraction of lines from the cache that had to be written back to main memory (in a copy-back cache) ranged from 17 to 56 percent.

Several issues bear on the trade-off between write-through and copy-back.

1. **Main Memory Traffic.** Copy-back almost always results in less main memory traffic since write-through requires a main memory access on every store, whereas copy-back only requires a store to main memory if the swapped out line (when a miss occurs) has been modified. Copy-back generally, though, results in the entire line being written back, rather than just one or two words, as would occur for each write memory reference (unless "dirty bits" are associated with partial lines; a dirty bit, when set, indicates the line has been modified while in the cache). For example, assume a miss ratio of 3 percent, a line size of 32 bytes, a memory module width of 8 bytes, a 16 percent store frequency, and 30 percent of all cache lines requiring a copy-back operation. Then the ratio of main memory store cycles to total memory references is 0.16 for write-through and 0.036 for copy-back.

2. **Cache Consistency.** If store-through is used, main memory always contains an up-to-date copy of all information in the system. When there are multiple processors in the system (including independent channels), main memory can serve as a common and consistent storage place, provided that additional mechanisms are used. Otherwise, either the cache must be shared or a complicated directory system must be employed to maintain consistency. This subject is discussed further in Section 2.7, but we note here that store-through simplifies the memory consistency problem.

3. **Complicated Logic.** Copy-back may complicate the cache logic. A dirty bit is required to determine when to copy a line back. In addition, arrangements have to be made to perform the copy-back before the fetch (on a miss) can be completed.

4. **Fetch-on-write.** Using either copy-back or write-through still leaves undecided the question of whether to fetch-on-write or not, if the information referenced is not in the cache. With copy-back, one will usually fetch-on-write, and with write-through, usually not. There are additional related possibilities and problems. For example, when using write-through, one could not only not fetch-on-write but one could choose actually to purge the modified line from the cache should it be found there. If the line is found in the cache, its replacement status (e.g., LRU) may or may not be updated. This is considered in item 6 below.

5. **Buffering.** Buffering is required for both copy-back and write-through. In copy-back, a buffer is required so that the line to be copied back can be held temporarily in order to avoid interfering with the fetch. One optimization worth noting for copy-back is to use spare cache/main memory cycles to do the copy-back of "dirty" (modified) lines [Baird81b]. For write-through, it is important to buffer several stores so that the CPU does not have to wait for them to be completed. Each buffer consists of a data part (the data to be stored) and the address part (the target address). In Smrt79 it was shown that a buffer with capacity of four provided most of the performance improvement possible in a write-through system. This is the number used in the IBM 3033. We note that a great deal of extra logic may be required if buffering is used. There is not only the logic required to implement the buffers, but also there must be logic to test all memory access addresses and match them against the addresses in the address part of the buffers. That is, there may be accesses to the material contained in the store buffers before the data in those buffers has been transferred to main memory. Checks must be made to avoid possible inconsistencies.

6. **Reliability.** If store-through is used, main memory always has a valid copy of the total memory state at any given time. Thus, if a processor fails (along with its cache), a store-through system can often be restored more easily. Also, if the only valid copy of a line is in the cache, an error-correcting code is needed there. If a cache error can be corrected from main memory, then a parity check is sufficient in the cache.

Some experiments were run to look at the miss ratio for store-through and copy-back. A typical example is shown in Figure 22; the other sets of traces yield very similar results. (In the case of write-through, we have counted each write as a miss.) It is clear that write-through always produces a much higher miss ratio. The terms *reorder* and *no reorder* specify how the replacement status of the lines were updated. *Reorder* means that a modified line is

WRITE POLICY EXPERIMENTS

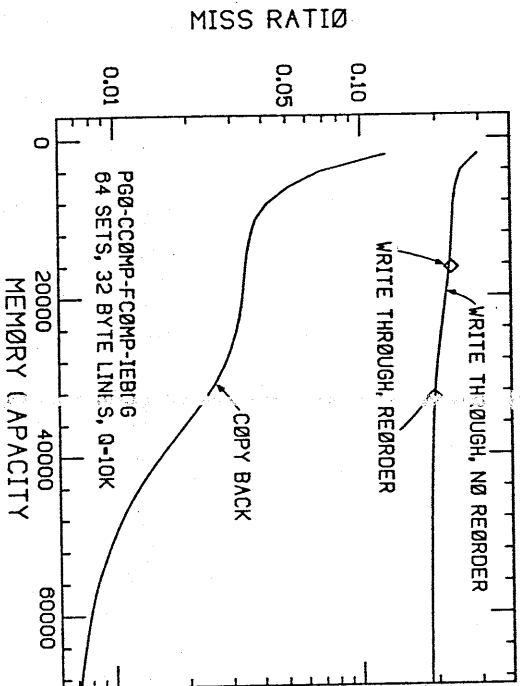


Figure 22. Miss ratio for copy-back and write-through. All writes counted as misses for write-through. No reordering implies replacement status of modified on write; reordering implies replacement status is updated on write.

moved to the top of the LRU stack within its set in the cache. No reordering implies that the replacement status of the line is not modified on write. From Figure 22, it can be seen that there is no significant difference in the two policies. For this reason, the IBM 3033, a write-through machine, does not update the LRU status of lines when a write occurs.

With respect to performance, there is no clear choice to be made between write-through and copy-back. This is because a good implementation of write-through selection has to wait for a write to complete. A good implementation of write-through rewrites, however, both sufficient buffering of writes and sufficient interleaving of main memory so that the probability of the CPU becoming blocked while waiting on a store is small. This appears to be true in the IBM 3033, but at the expense of a great deal of buffering and other logic. For example, in the 3033 each buffered store requires a double-word datum buffer, a single buffer for the store address, and a 1-byte buffer to

indicate which bytes have been modified in the double-word store. There are also components to match each store address against subsequent accesses to main memory so that references to the modified data get the updated values. Copy-back could probably have been implemented much more cheaply.

The Andahl Computers all use copy-back, as does the IBM 3081. IBM uses store-through in the 370/168 [IBM75] and the 3033 [IBM78], as does DEC in the PDP-11/70 [Stre76] and VAX 11/780 [DEC78]. Honeywell in the 66/60 and 66/80, and Icl in the AS/6.

2.6 Effect of Multiprogramming: Cold-Start and Warm-Start

A significant factor in the cache miss ratio is the frequency of task switching, or inversely, the value of the mean intertask switching time, Q . The problem with task switching is that every time the active task is changed, a new process may have to be

loaded from scratch into the cache. This issue was discussed by East75, where the terms *warm-start* and *cold-start* were coined to refer to the miss ratio starting with a full memory and the miss ratio starting with an empty memory, respectively. Other papers which discuss the problem include East78, Kona80, MacD79, Peur77, Ponn75, Schm71, and Stre76.

Typically, a program executes for a period of time before an interruption (I/O, clock, etc.) of some type invokes the supervisor. The supervisor eventually relinquishes control of the processor to some user process, perhaps the same one as was running most recently. If it is not the same user process, the new process probably does not find any lines of its address space in the cache, and starts immediately with a number of misses. If the most recently executed process is restarted, and if the supervisor interruption has not been too long, some useful information may still remain. In Peur77, some figures are given about the length of certain IBM operating system supervisor interruptions and what fraction of the cache is purged. (One may also view the user as interrupting the supervisor and increasing the supervisor miss ratio.)

The effect of the task-switch interval on the miss ratio cannot be easily estimated. In particular, the effect depends on the workload and on the cache size. We also observe that the proportion of cache misses due to task switching increases with increasing cache size, even though the absolute miss ratio declines. This is because a small cache has a large inherent miss ratio (since it does not hold the program's working set) and this miss ratio is only slightly augmented by task-switch-induced misses. Conversely, the inherent low miss ratio of a large cache is greatly increased, in relative terms, by task switching. We are not aware of any current machine for which a breakdown of the miss ratio into these two components has been done.

Some experimental results bearing on this problem appear in Figures 23 and 24. In each, the miss ratio is shown as a function of the memory size and task-switch interval Q (Q is the number of memory references). The figures presented can be understood as follows. A very small Q (e.g., 100, 1,000) implies that the cache is shared

between all of the active processes, and that when a process is restarted it finds a significant fraction of its previous information still in the cache. A very large Q (e.g., 100,000, 250,000) implies that when the program is restarted it finds an empty cache (with respect to its own working set), but that the new task runs long enough first to fill the cache and then to take advantage of the full cache. Intermediate values for Q result in the situation where a process runs for a while but does not fill the cache; however, when it is restarted, none of its information is still cache resident (since the multiprogramming degree is four). These three regions of operation are evident in Figures 23 and 24 as a function of Q and of the cache size. (In Sarr78, Q is estimated to be about 25,000.)

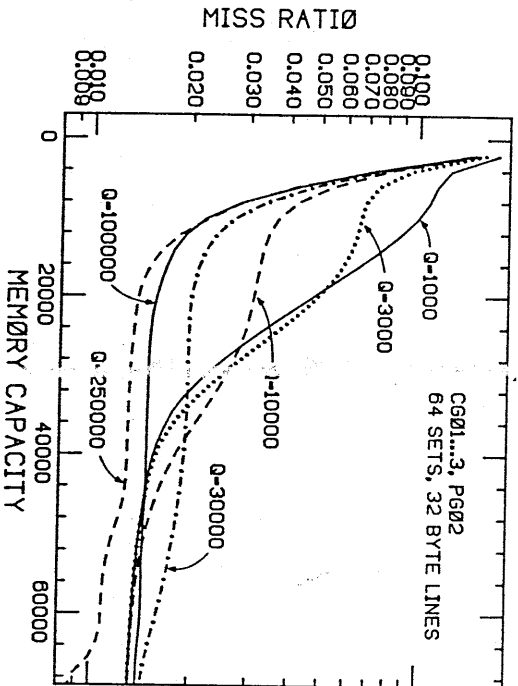
There appear to be several possible solutions to the problem of high cache miss ratios due to task switching. (1) It may be possible to lengthen the task-switch interval. (2) The cache can be made so large that several programs can maintain information in it simultaneously. (3) The scheduling algorithm may be modified in order to give preference to a task likely to have information resident in the cache. (4) If the working set of a process can be identified (e.g., from the previous period of execution), it might be reloaded as a whole; this is called *working set restoration*. (5) Multiple caches may be created; for example, a separate cache could be established for the supervisor to use, so that, when invoked, it would not displace user data from the cache. This idea is considered in Section 2.10, and some of the problems of the approach are indicated.

A related concept is the idea of bypassing the cache for operations unlikely to result in the reuse of data. For example, long vector operations such as a very long move character (e.g., IBM MYCCL) could bypass the cache entirely [Losq82] and thereby avoid displacing other data more likely to be reused.

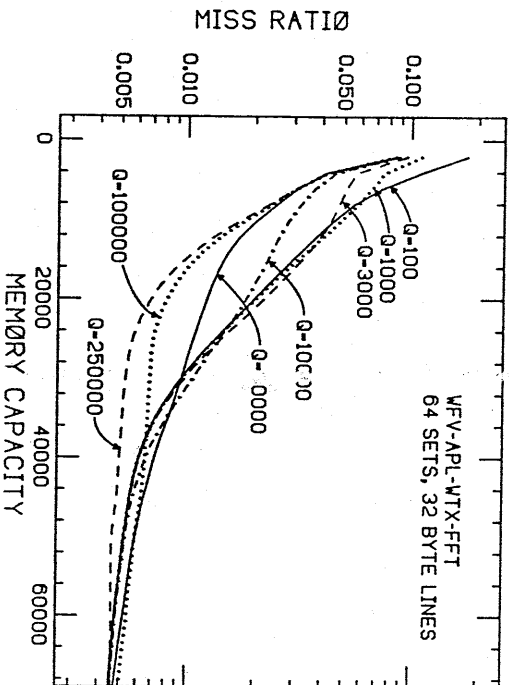
2.7 Multicache Consistency

Large modern computer systems often have several independent processors, consisting sometimes of several CPUs and sometimes of just a single CPU and several channels.

VARY MP QUANTUM SIZE

Figure 23. Miss ratio versus memory capacity or range of multiprogramming intervals Q .

VARY MP QUANTUM SIZE

Figure 24. Miss ratio versus memory capacity or range of multiprogramming intervals Q .

Each processor may have zero, one, or several caches. Unfortunately, in such a multiple processor system, a given piece of information may exist in several places at a given time, and it is important that all processors have access (as necessary) to the same, unique (at a given time) value. Several solutions exist and/or have been proposed for this problem. In this section, we discuss many of these solutions; the interested reader should refer to BEAN⁷⁹, CENS⁷⁸, DRUM^{81a}, DUBO⁸², JONE⁷⁶, JONE^{77a}, MAZA⁷⁷, MCW⁷⁷, NGA⁸¹, and TANG⁷⁶ for additional explanation.

As a basis for discussion, consider a single CPU with a cache and with a main memory behind the cache. The CPU reads an item into the cache and then modifies it. A second CPU (of similar design and using the same main memory) also reads the item and modifies it. Even if the CPUs were using store-through, the modification performed by the second CPU would not be reflected in the cache of the first CPU unless special steps were taken. There are several possible special steps.

1. Shared Cache. All processors in the system can use the same cache. In general, this solution is infeasible because the bandwidth of a single cache usually is not sufficient to support more than one CPU, and because additional access time delays may be incurred because the cache may not be physically close enough to both (or all) processors. This solution is employed successfully in the Amdahl 470 computers, where the CPU and the channels all use the same cache; the 470 series does not, however, permit tightly coupled CPUs. The UNIVAC 1100/80 [BORC⁷⁹] permits two CPUs to share one cache.

2. Broadcast Writes. Every time a CPU performs a write to the cache, it also sends that write to all other caches in the system. If the target of the store is found in some other cache, it may be either updated or invalidated. Invalidation may be less likely to create inconsistencies, since updates can possibly "cross" such that CPU A updates its own cache and then B's cache. CPU B simultaneously updates its own cache and then A's. Updates also require more data transfer. The IBM 370/168 and 3083 processors use invalidation. A store by a CPU

or channel is broadcast to all caches sharing the same main memory. This broadcast store is placed in the *buffer invalidation address stack* (BIAS) which is a list of addresses to be invalidated in the cache. The buffer invalidation address stack has a high priority for cache cycles, and if the target line is found in that cache, it is invalidated.

The major difficulty with broadcasting store addresses is that every cache memory in the system is forced to surrender a cycle for invalidation lookup to any processor which performs a write. The memory interference that occurs is generally acceptable for two processors (e.g., IBM's current MP systems), but significant performance degradation is likely with more than two processors. A clever way to minimize this problem appears in a recent patent [BEAN⁷⁹]. In that patent, a BIAS Filter Memory (BFM) is proposed. A BFM is associated with each cache in a tightly coupled MP system. This filter memory works by filtering out repeated requests to invalidate the same block in a cache.

3. Software Control. If a system is being written from scratch and the architecture can be designed to support it, then software control can be used to guarantee consistency. Specifically, certain information can be designated noncacheable, and can be accessed only from main memory. Such items are usually semaphores and perhaps data structures such as the job queue. For efficiency, some shared writable data has to be cached. The CPU must therefore be equipped with commands that permit it to purge any such information from its own cache as necessary. Access to shared writable cacheable data is possible only within critical sections, protected by noncacheable semaphores. Within the critical regions, the code is responsible for restoring all modified items to main memory before releasing the lock. Just such a scheme is intended for the S-1 multiprocessor system under construction at the Lawrence Livermore Laboratory [HAU⁷⁹, MCW⁷⁷]. The Honeywell Series 66 machines use a similar mechanism. In some cases, the simpler alternative of making shared writable information noncacheable may be acceptable.

4. Directory Methods. It is possible to keep a centralized and/or distributed direc-

tory of all main memory lines, and use it to ensure that no lines are write-shared. One such scheme is as follows, though several variants are possible. The main memory maintains $k + 1$ bits for each line in main memory, when there are k caches in the system. Bit i , $i = 1, \dots, k$ is set to 1 if the corresponding cache contains the line. The $(k + 1)$ th bit is 1 if the line is being or has been modified in a cache, and otherwise is 0. If the $(k + 1)$ th bit is on, then exactly one of the other bits is on. Each CPU has associated with each line in its cache a single bit (called the *private bit*). If that bit is on, that CPU has the only valid copy of that line. If the bit is off, other caches and main memory may also contain current copies. Exactly one private bit is set if and only if the main memory directory bit $k + 1$ is set.

A CPU can do several things which provoke activity in this directory system. If a CPU attempts to read a line which is not in its cache, the main memory directory is queried. There are two possibilities: either but $k + 1$ is off, in which case the line is transferred to the requesting cache and the corresponding bit set to indicate this; or, bit $k + 1$ is on, in which case the main memory directory must recall the line from the cache which contains the modified copy, update main memory, invalidate the copy in the cache that modified it, send the line to the requesting CPU/cache and finally update itself to reflect these changes. (Bit $k + 1$ is then set to zero, since the request was a read.)

An attempt to perform a write causes one of three possible actions. If the line is already in cache and has already been modified, the private bit is on and the write takes place immediately. If the line is not in cache, then the main memory directory must be queried. If the line is in any other cache, it must be invalidated (in all other caches), and main memory must be updated if necessary. The main memory directory is then set to reflect the fact that the new cache contains the modified copy of the data, the line is transmitted to the requesting cache, and the private bit is set. The third possibility is that the cache already contains the line but that it does not have its private bit set. In this case, permission must be requested from the main

memory directory to perform the write. The main memory directory invalidates any other copies of the line in the system, marks its own directory suitably, and then gives permission to modify the data.

The performance implications of this method are as follows. The cost of a miss may increase significantly due to the need to query the main memory directory and possibly retrieve data from other caches. The use of shared writable information becomes expensive due to the high miss ratio: that are likely to be associated with such information. In CENS78, there is some attempt at a quantitative analysis of these performance problems.

Another problem is that I/O overruns may occur. Specifically, an I/O data stream may be delayed while directory operations take place. In the meantime, some I/O data are lost. Care must be taken to avoid this problem. Either substantial I/O buffering or write-through is clearly needed.

Other variants of method 4 are possible. (1) The purpose of the central directory is to minimize the queries to the caches of other CPUs. The central directory is not logically necessary; sufficient information exists in the individual caches. It is also possible to transmit information from cache to cache without going through main memory. (2) If the number of main memory directory bits is felt to be too high, locking could be on a page instead of on a line basis. (3) Fore-through may be used instead of copy back; thus main memory always has a valid copy and data do not have to be fetched from the other caches, but can simply be invalidated in those other caches.

The IBM 3081D, which contains two CPUs, essentially uses the directory scheme described. The higher performance 3081K functions similarly, but passes the necessary information from cache to cache rather than going through main memory.

Another version of the directory method is called the broadcast search [DRM81b]. In it is case, a miss is sent not only to the main memory but to all caches. Whichever caches (cache or main) contain the desired information send it to the requesting processor.

Lin [Lin82] proposes a multicache scheme to minimize the overhead of directory operations. He suggests that all CPUs

have two caches, only one of which can contain shared data. The overhead of directory access and maintenance would thus only be incurred when the shared data cache is accessed.

There are two practical methods among the above alternatives: method 4 and the BIAS Filter Memory version of method 2. Method 4 is quite general, but is potentially very complex. It may also have performance problems. No detailed comparison exists, and other and better designs may yet remain to be discovered. For new machines, it is not known whether software control is better than hardware control; clearly, for existing architectures and software, hardware control is required.

2.8 Data/Instruction Cache

Two aspects of the cache having to do with its performance are cache bandwidth and access time. Both of these can be improved by splitting the cache into two parts, one for data and one for instructions. This doubles the bandwidth since the cache can now service two requests in the time it formerly required for one. In addition, the two requests served are generally complementary. Fast computers are pipelined, which means that several instructions are simultaneously in the process of being decoded and executed. Typically, there are several stages in a pipeline, including instruction fetch, instruction decode, operand address generation, operand fetch, execution, and transmission of the results to their destination (e.g., to a register). Therefore, while one instruction is being fetched (from the instruction cache), another can be having its operands fetched from the operand cache. In addition, the logic that arbitrates priority between instruction and data accesses to the cache can be simplified or eliminated.

A split instruction/data cache also provides access time advantages. The CPU of a high-speed computer typically contains (exclusive of the S-unit) more than 100,000 logic gates and is physically large. Further, the logic having to do with instruction fetch and decode has little to do with operand fetch and store except for *execute* instructions and possibly for the targets of branches. With a single cache system, it is not always possible simultaneously to place the cache immediately adjacent to all of the

logic which will access it. A split cache, on the other hand, can have each of its halves placed in the physical location which is most useful, thereby saving from a fraction of a nanosecond to several nanoseconds.

There are, of course, some problems introduced by the split cache organization. First, there is the problem of consistency. Two copies now exist of information which formerly existed only in one place. Specifically, instructions can be modified, and this modification must be reflected before the instructions are executed. Further, it is possible that even if the programs are not self-modifying, both data and instructions may coexist in the same line, either for reasons of storage efficiency or because of immediate operands. The solutions for this problem are the same as those discussed in the section on multicache consistency (Section 2.7), and they work here as well. It is imperative that they be implemented in such a way so as to not impair the access time advantage given by this organization.

Another problem of this cache organization is that it results in inefficient use of the cache memory. The size of the working set of a program varies constantly, and in particular, the fraction devoted to data and to instructions also varies. (A dynamically split design is suggested by FAYR78.) If the instructions and data are not stored together, they must each exist within their own memory, and be unable to share a larger amount of that resource. The extent of this problem has been studied both experimentally and analytically. In SHED76, a set of formulas are provided which can be used to estimate the performance of the unified cache from the performance of the individual ones. The experimental results were not found to agree with the mathematical ones, although the reason was not investigated. We believe that the nonstationarity of the workload was the major problem.

We compared the miss ratio of the split cache to that of the unified cache for each of the sets of traces; some of the results appear in Figures 25-28. (See Bell74 and TRAK78 for additional results.) We note that there are several possible ways to split and manage the cache and the various alternatives have been explored. One could split the cache in two equal parts (labeled "SPLIT EQUAL"), or the observed miss

I/D CACHES VS. UNIFIED CACHE

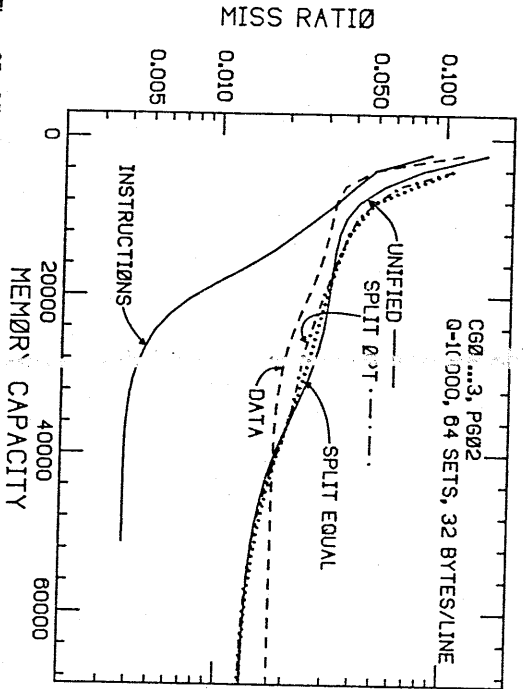


Figure 25. Miss ratio versus memory capacity for a unified cache, cache split equally between instruction and data halves, cache split according to a static optimum partition between instruction and data halves, and miss ratios individually for instruction and data halves.

I/D CACHES VS. UNIFIED CACHE

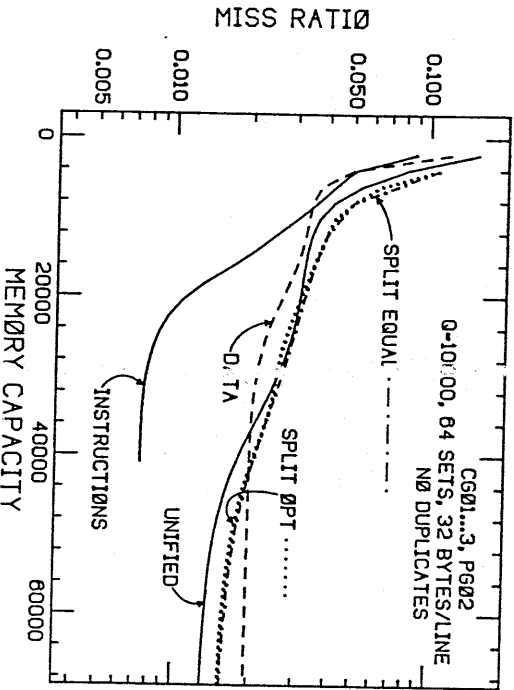


Figure 26. Miss ratio versus memory capacity for a unified cache, cache split equally between instruction and data halves, cache split according to a static optimum partition between instruction and data halves, and miss ratios individually for instruction and data halves.

I/D CACHES VS. UNIFIED CACHE

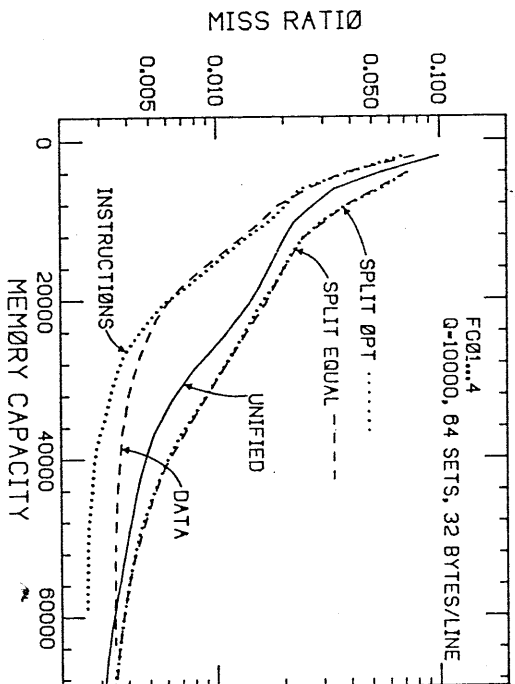


Figure 27. Miss ratio versus memory capacity for a unified cache, cache split equally between instruction and data halves, cache split according to a static optimum partition between instruction and data halves, and miss ratios individually for instruction and data halves.

I/D CACHES VS. UNIFIED CACHE

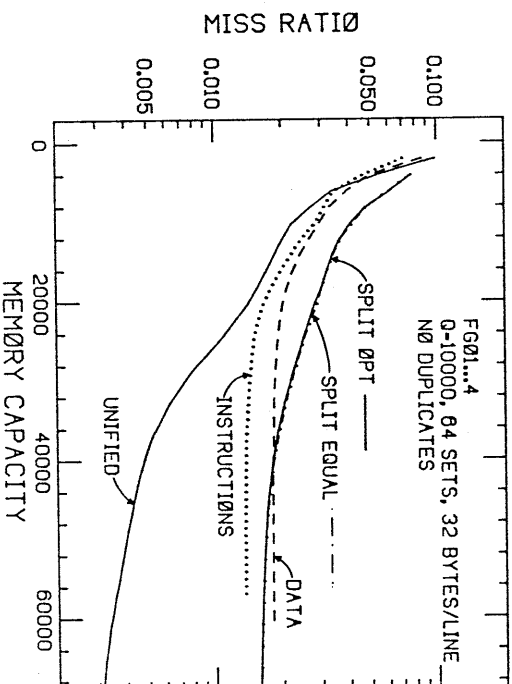


Figure 28. Miss ratio versus memory capacity for a unified cache, cache split equally between instruction and data halves, cache split according to a static optimum partition between instruction and data halves, and miss ratios individually for instruction and data halves.

ratios could be used (from this particular run) to determine the optimal static split ("SPLIT OPT"). Also, when a line is found to be in the side of the cache other than the one currently accessed, it could either be duplicated in the remaining side of the cache or it could be moved; this latter case is labeled "NO DUPLICATES." (No special label is shown if duplicates are permitted.) The results bearing on these distinctions are given in Figures 25-28.

We observe in Figures 25 and 26 that the unified cache, the equally split cache, and the cache split unequally (split optimally) all perform about equally well with respect to miss ratio. Note that the miss ratio for the instruction and data halves of the cache individually are also shown. Further, comparing Figures 25 and 26, it seems that barring duplicate lines has only a small negative effect on the miss ratio.

In sharp contrast to the measurements of Figures 25 and 26 are those of Figures 27 and 28. In Figure 27, although the equally and optimally split cache are comparable, the unified cache is significantly better. The unified cache is better by an order of magnitude when duplicates are not permitted (Figure 28), because the miss ratio is sharply increased by the constant movement of lines between the two halves. It appears that lines sharing instruction and data are very common in programs compiled with IBM's FORTRAN G compiler and are not common in programs compiled using the IBM COBOL or PL/I compiler. (Results similar to the FORTRAN case have been found for the two other sets of IBM traces, all of which include FORTRAN code but are not shown.)

Based on these experimental results, we can say that the miss ratio may increase significantly if the caches are split, but that the effect depends greatly on the workload. Presumably, compilers can be designed to minimize this effect by ensuring that data and instructions are in separate lines, and perhaps even in separate pages.

Despite the possible miss ratio penalty of splitting the cache, there are at least two experimental machines and two commercial ones which do so. The S-1 [Halt79, McWitt77] at Lawrence Livermore Laboratory is being built with just such a cache; it relies on (new) software to minimize the

problems discussed here. The 801 minicomputer, built at IBM Research (Yorktown Heights) [Elec76, Rad82] also has a split cache. The Hitachi H200 and Intel AS/6 [Ross79] both have a split data/instruction cache. No measurements have been publicly reported for any of these machines.

2.9 Virtual Address Cache

Most cache memories address the cache using the real address (see Figure 2). As the reader recalls, we discussed (Introduction, Section 2.3) the fact that the virtual address was translated by the TLB to the real address, and that the line lookup and readout could not be completed until the real address was available. This suggests that the cache access time could be significantly reduced if the translation step could be eliminated. The way to do this is to address the cache directly with the virtual address. We call a cache organized this way a virtual address cache. The MU-5 [Bane77] uses this organization for its name store. The S-1, the IBM 801, and the ICL 2900 series machines also use this idea. It is discussed in Eze79. See also Olab79.

There are some additional considerations in building a virtual address cache, and there is one serious problem. First, all addresses must be tagged with the identifier of the address space with which they are associated, or else the cache must be purged every time task switching occurs. Tagging is not a problem, but the address tag in the cache must be extended to include the address space ID. Second, the translation mechanism must still exist and must still be efficient, since virtual addresses must be translated to real addresses whenever main memory is accessed, specifically for misses and for writes in a write-through cache. The TLB cannot be eliminated.

The most serious problem is that of "synonyms," two or more virtual addresses that map to the same real address. Synonyms occur whenever two address spaces share a code or data. (Since the lines have address space tags, the virtual addresses are different even if the line occurs in the same place in both address spaces.) Also, the supervisor may exist in the address space of each user, and it is important that

only one copy of supervisor tables be kept. The only way to detect synonyms is to take the virtual address, map it into a real address, and then see if any other virtual addresses in the cache map into the same real address. For this to be feasible, the inverse mapping must be available for every line in the cache; this inverse mapping is accessed on real address and indicates all virtual addresses associated with that real address. Since this inverse mapping is the opposite of the TLB, we choose to call the inverse mapping buffer (if a separate one is used) the RTB or reverse translation buffer. When a miss occurs, the virtual address is translated into the real address by the TLB. The access to main memory for the miss is overlapped with a similar search of the RTB to see if that line is already in the cache under a different name (virtual address). If it is, it must be renamed and moved to its new location, since multiple copies of the same line in the cache are clearly undesirable for reasons of consistency.

The severity of the synonym problem can be decreased if shared information can be forced to have the same location in all address spaces. Such information can be given a unique address space identifier, and the lookup algorithm always considers such a tag to match the current address space ID. A scheme like this is feasible only for the supervisor since other shared code should not conceivably be so allocated. Shared supervisor code does have a unique location in all address spaces in IBM's MVS operating system.

The RTB may or may not be a simple structure, depending on the structure of the rest of the cache. In one case it is fairly simple: if the bits used to select the set of the cache are the same for the real and virtual address (i.e., if none of the bits used to select the set undergo translation), the RTB can be implemented by associating with each cache line two address tags [Bane79]. If a match is not found on the virtual address, then a search is made in the tag set on the real address. If that search finds the line, then the virtual address tag is changed to the current virtual address. A more complex design would involve a separate mapping buffer for the reverse translation.

2.10 User/Supervisor Cache

It was suggested earlier that a significant fraction of the miss ratio is due to task-switching. A possible solution to this problem is to use a cache which has been split into two parts, one of which is used only by the supervisor and the other of which is used primarily by the user state programs. If the scheduler were programmed to restart, when possible, the same user program that was running before an interrupt, then the user state miss ratio would drop appreciably. Further, if the same interrupts recur frequently, the supervisor state miss ratio may also drop. In particular, neither the user nor the supervisor would purge the cache of the other's lines. (See Peur77 for some data relevant to this problem.) The supervisor cache may have a high miss ratio in any case due to its large working set. (See Mila75 for an example.)

Despite the appealing rationale of the above comments, there are a number of problems with the user/supervisor cache. First, it is actually unlikely to cut down the miss ratio. Most misses occur in supervisor state [Mila75] and a supervisor cache half as large as the unified cache is likely to be worse since the maximum cache capacity is no longer available to the supervisor. Further, it is not clear what fraction of the time the scheduling algorithm can restart the same program. Second, the information used by the user and the supervisor are not entirely distinct, and cross-access must be permitted. This overlap introduces the problem of consistency.

We are aware of only one evaluation of the split user/supervisor cache [Ross80]. In that case, an experiment was run on an Hitachi M180. The results seemed to show that the split cache performed about as well as a unified one, but poor experimental design makes the results questionable. We do not expect that a split cache will prove to be useful.

2.11 Input/Output through the Cache

In Section 2.7, the problem of multicache consistency was discussed. We noted that if all accesses to main memory use the same cache, then there would be no consistency problem. Precisely this approach has been

used in one manufacturer's computers (Amdahl Corporation).

2.11.1 Overruns

While putting all input/output through the cache solves the consistency problem, it introduces other difficulties. First, there is the overrun problem. An overrun occurs when for some reason the I/O stream cannot be properly transmitted between the memory (cache or main) and the I/O device. Transmitting the I/O through the cache can cause an overrun when the line accessed by the I/O stream is not in the cache (and is thus a miss) and for some reason cannot be obtained quickly enough. Most I/O devices involve physical movement, and when the buffering capacity embedded in the I/O path is exhausted, the transfer must be aborted and then restarted. Overruns can be provoked when:

- (1) the cache is already processing one or more misses and cannot process the current (additional) one quickly enough;
- (2) more than one I/O transfer is in progress, and more active (in use) lines map into one set than the set size can accommodate; or
- (3) the cache bandwidth is not adequate to handle the current burst of simultaneous I/O from several devices. Overruns can be minimized if the set size of the cache is large enough, the bandwidth is high enough, and the ability to process misses is sufficient. Sufficient buffering should also be provided in the I/O paths to the devices.

2.11.2 Miss Ratio

Directing the input/output data streams through the cache also has an effect on the miss ratio. This I/O data occupies some fraction of the space in the cache, and this increases the miss ratio for the other users of the cache. Some experiments along these lines were run by the author and results are shown in Figures 29-32. IORATE is the ratio of the rate of I/O accesses to the cache to the rate of CPU accesses. (I/O activity is simulated by a purely sequential synthetic address stream referencing a distinct address space from the other programs.) The miss ratio as a function of memory size and I/O transfer rate is shown in Figures 29 and 30 for two of the sets of traces. The

data has been rearranged to show more directly the effect on the miss ratio in Figures 29 and 32. The results displayed here show no clear mathematical pattern, and we were unable to derive a useful and verifiable formula to predict the effect on the miss ratio by an I/O stream.

Examination of the results presented in Figures 29-32 suggests that for reasonable I/O rates (less than 0.05; see POWER7 for some I/O rate data) the miss ratio is not affected to any large extent. This observation is consistent with the known performance of the Amdahl computers, which are not seriously degraded by high I/O rates.

2.12 Cache Size

Two very important questions when selecting a cache design are how large should the cache be and what kind of performance can we expect. The cache size is usually dictated by a number of criteria having to do with the cost and performance of the machine. The cache should not be so large that it represents an expense out of proportion to the added performance, nor should it occupy an unreasonable fraction of the physical space within the processor. A very large cache may also require more access circuitry, which may increase access time.

As a result of the warnings given in the paragraph above, one can generally assume that the larger the cache, the higher the hit ratio and therefore the better the performance. The issue is then one of the relation between cache size and hit ratio. This is a very difficult problem, since the cache hit ratio varies with the workload and the machine architecture. A cache that might yield a 99 percent hit ratio on a PDP-11 program could result in a 90 percent or lower hit ratio for IBM (MVS) supervisor state code. This problem cannot be usefully studied using trace-driven simulation because the miss ratio varies tremendously from program to program and only a small number of traces can possibly be analyzed. Typical trace-driven simulation results appear throughout this paper, however, and the reader may wish to scan that data for insight. There is also a variety of data available in the literature and the reader may wish to inspect the results presented in Als*78, Bell*74, Berg*6, Gibs*7, Lee*69,

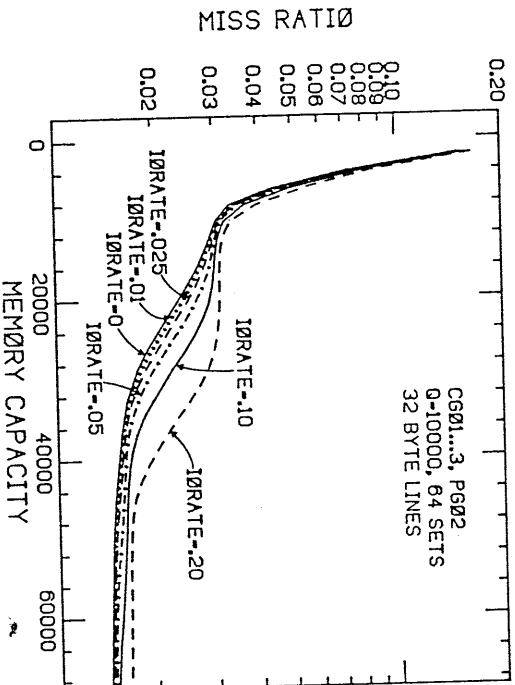


Figure 29. Miss ratio versus memory capacity while I/O occurs through cache at specified rate. IORATE refers to fraction of all memory references due to I/O.

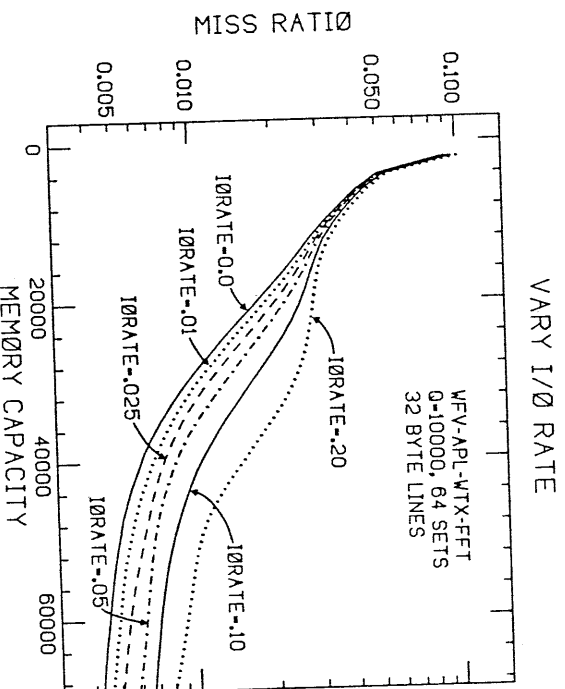


Figure 30. Miss ratio versus memory capacity while I/O occurs through cache.

MISS RATIO VS. I/O RATE

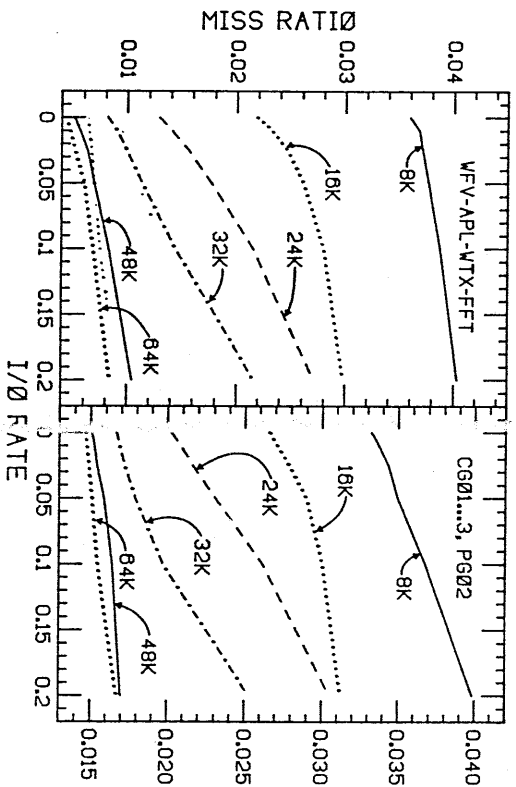


Figure 31. Miss ratio versus I/O rate for variety of memory capacities.

INCREASE IN MISS RATIO VS. I/O RATE

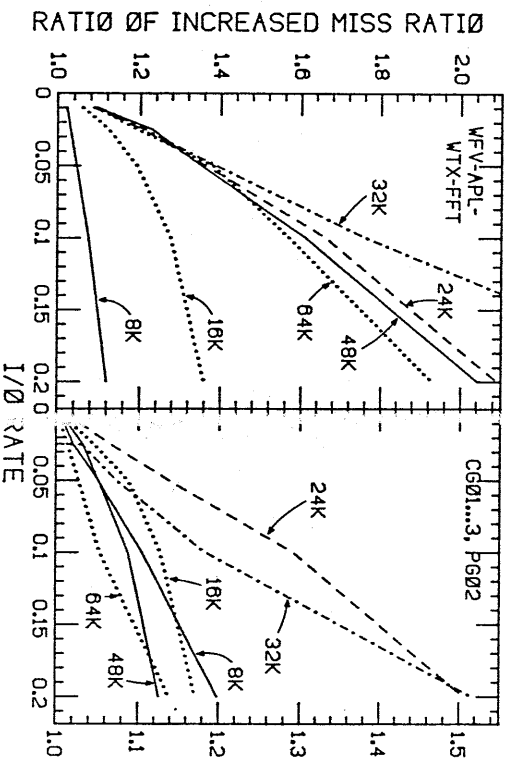


Figure 32. Miss ratio versus I/O rate for variety of memory capacities.

MISS RATIO VS. CACHE SIZE

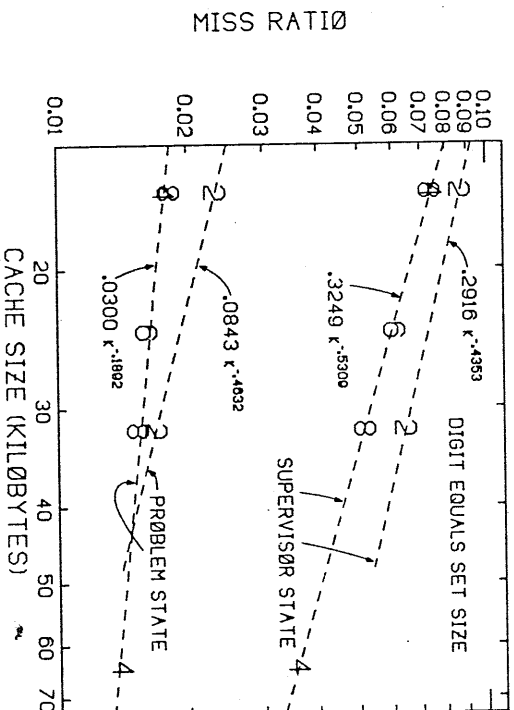


Figure 33. Miss ratio versus cache size, classified by set size and user/supervisor state. Data gathered by hardware monitor from machine running standard benchmark.

MEAD70, STRE76, THAK78, and YUVA75. Some direct measurements of cache miss ratios appear in CLAR81 and CLAR82. In the former, the Dorado was found to have hit ratios above 99 percent. In the latter, the VAX 11/780 was found to have hit ratios around 90 percent.

Another problem with trace-driven simulation is that in general user state programs are the only ones for which many traces exist. In IBM MVS systems, the supervisor typically uses 25 to 60 percent of the CPU time, and provides by far the largest component of the miss ratio [MILA75]. User programs generally have very low miss ratios, in our experience, and many of those misses come from task switching.

Two models have been proposed in the literature for memory hierarchy miss ratios. Saltzer [SALT74] suggested, based on his data, that the mean time between faults was linearly related to the capacity of the memory considered. But later results, taken on the same system [GREF74] contradict

Saltzer's earlier findings. [CHOW75 and CHOW76] suggest that the miss ratio curve was of the form $m = a(c^b)$, where a and b are constants, m is the miss ratio, and c is the memory capacity. They show no experimental results to substantiate this model, and it seems to have been chosen for mathematical convenience.

Actual cache miss ratios, from real machines running "typical" workloads, are the most useful source of good measurement data. In Figure 33 we show a set of such measurements taken from Amdahl 470 computers running a standard Amdahl internal benchmark. This data is reproduced from HARD80a. Each digit represents a measurement point, and shows either the supervisor or problem state miss ratio for a specific cache size and set size; the value of the digit at each point is the set size. Examination of the form of the measurements from Figure 33 suggest that the miss ratio can be approximated over the range shown by an equation of the form $m = a(c^b)$ (consistent with CHOW75 and CHOW76),

where m is the miss ratio, α and b are constants ($b < 0$), and k is the cache capacity in kilobytes. The values of α and b are shown for four cases in Figure 33; supervisor and user state for a set size of two, and supervisor and user state for all other set sizes. We make no claims for the validity of this function for other workloads and/or other architectures, nor for cache sizes beyond the range shown. From Figure 33 it is evident, though, that the supervisor contributes the largest fraction of the miss ratio, and the supervisor state measurements are quite consistent. Within the range indicated, therefore, these figures can probably be used for a first approximation at estimating the performance of a cache design.

Typical cache sizes in use include 128 kbytes (NEC ACOS 9000), 64 kbytes (Amdahl 470V/8, IBM 3083, IBM 3081K per CPU), 32 kbytes (IBM 370/168-3, IBM 3081D per CPU, Amdahl 470V/7, Magnuson M80/43), 16 kbytes (Amdahl 470V/6, Intel AS/6, IBM 4341, Magnuson M80/42, M80/44, DEC VAX 11/780), 8 kbytes (Honeywell 66/60 and 66/80, Xerox Dorado [C1A81]), 4 kbytes (VAX 11/750, IBM 4331), 1 kbyte (PDP-11/70).

2.13. Cache Bandwidth, Data Path Width, and Access Resolution

2.13.1 Bandwidth

For adequate performance, the cache bandwidth must be sufficient. Bandwidth refers to the aggregate data transfer rate, and is equal to data path width divided by access time. The bandwidth is important as well as the access time, because (1) there may be several sources of requests for cache access (instruction fetch, operand fetch and store, channel activity, etc.), and (2) some requests may be for a large number of bytes. If there are other sources of requests for cache cycles, such as prefetch lookups and transfers, it must be possible to accommodate these as well.

In determining what constitutes an adequate data transfer rate, it is not sufficient that the cache bandwidth exceed the average demand placed on it by a small amount. It is important as well to avoid contention since if the cache is busy for a cycle and one or more requests are blocked, these blocked

requests can result in permanently wasted machine cycles. In the Amdahl 470V/8 and the IBM 3083, the cache bandwidth appears to exceed the average data rate by a factor of two to three, which is probably the minimum sufficient margin. We note that in the 470V/6 (when prefetch is used experimentally) prefetches are executed only during otherwise idle cycles, and it has been observed that not all of the prefetches actually are performed. (Newly arrived prefetch requests take the place of previously queued but never performed requests.)

For some instructions, the cache bandwidth can be extremely important. This is particularly the case for data movement instructions such as: (1) instructions which load or unload all of the registers (e.g., IBM 370: instructions STM, LM); (2) instructions which move long character strings (MVC, MV, IL); and (3) instructions which operate on long character strings (e.g., CLC, OC, NC, XC). In these cases, especially the first two, there is little if any processing to be done; the question is simply one of physical data movement, and it is important that the cache data path be as wide as possible—in large machines, 8 bytes (3083, 3081, Intel AS/6) instead of 4 (470V/6, V/7, V/8); in small machines, 4 bytes (VAX 11/780) instead of 2 (PDP-11/70).

It is important to note that cache data path width is expensive. Doubling the path width means doubling the number of lines into and out of the cache (i.e., the bus width) and all of the associated circuitry. This frequently implies some small increase in access time, due to larger physical packaging and/or additional levels of gate delay. The store, both the cost and performance aspects of cache bandwidth must be considered during the design process.

Another approach to increasing cache bandwidth is to interleave the cache [Dri 80, Yam80]. If the cache is required to serve a large number of small requests very quickly, it may be efficient to replicate the cache (e.g., two or four times) and access each separately, depending on the lower order bits of the desired locations. This approach is very expensive, and to our knowledge, has not been used on any existing machine. (See Pom75 for some additional comments.)

2.13.2 Priority Arbitration

An issue related to cache bandwidth is what to do when the cache has several requests competing for cache cycles and only one can be served at a given time. There are two criteria for making the choice: (1) give priority to any request that is "deadline scheduled" (e.g., an I/O request that would otherwise abort); and (2) give priority (after (1)) to requests in order to enhance machine performance. The second criterion may be sacrificed for implementation convenience, since the optimal scheduling of cache accesses may introduce unreasonable complexity into the cache design. Typically, fixed priorities are assigned to competing cache requests, but dynamic scheduling, though complex, is possible [Blou80].

An illustration of cache priority resolution, we consider two large, high-speed computers: the Amdahl 470V/7 and the IBM 3083. In the Amdahl machine, there are five "ports" or address registers, which hold the addresses for cache requests. Thus, there can be up to five requests queued for access. These ports are the operand port, the instruction port, the channel port, the translate port, and the prefetch port. The first three are used respectively for operand store and fetch, instruction fetch, and channel I/O (since channels use the cache also).

The translate port is used in conjunction with the TLB and translator to perform virtual to real address translation. The prefetch port is for a number of special functions, such as setting the storage key or purging the TLB, and for prefetch operations. There are sixteen priorities for the 470V/6 cache; we list the important ones here in decreasing order of access priority: (1) move line in from main storage, (2) operand store, (3) channel store, (4) fetch second half of double word request, (5) move line out from cache to main memory, (6) translate, (7) channel fetch, (8) operand fetch, (9) instruction fetch, (10) prefetch.

The IBM 3083 has a similar list of cache access priorities [IBM78]: (1) main memory fetch transfer, (2) invalidate line in cache modified by channel or other CPU, (3) search for line modified by channel or other CPU, (4) buffer reset, (5) translate, (6) redo

(some cache accesses are blocked and have to be restarted), and (7) normal instruction or operand access.

2.14 Multilevel Cache

The largest existing caches (to our knowledge) can be found in the NEC ACOS 9000 (128 kbytes) and the Amdahl 470V/8 and IBM 3083 processors (64 kbytes). Such large caches pose two problems: (1) their physical size and logical complexity increase the access time, and (2) they are very expensive. The cost of the chips in the cache can be a significant fraction (5-20 percent) of the parts cost of the CPU. The reason for the large cache, though, is to decrease the miss ratio. A possible solution to this problem is to build a two-level cache, in which the smaller, faster level is on the order of 4 kbytes and the larger, slower level is on the order of 64-512 kbytes. In this way, misses from the small cache could be satisfied, not in the six to twelve machine cycles commonly required, but in two to four cycles. Although the miss ratio from the small cache would be fairly high, the improved cycle time and decreased miss penalty would yield an overall improvement in performance. Suggestions to this effect may be found in Benn82, Ohno77, and Spar78. It has also been suggested for the TLB [Nca82].

As might be expected, the two-level or multilevel cache is not necessarily desirable. We suggested above that misses from the fast cache to the slow cache could be serviced quickly, but detailed engineering studies are required to determine if this is possible. The five-to-one or ten-to-one ratio of main memory to cache memory access times is not wide enough to allow another level to be easily placed between them.

Expense is another consideration. A two-level cache implies another level of access circuitry, with all of the attendant complications. Also, the large amount of storage in the second level, while cheaper per bit than the low-level cache, is not inexpensive on the whole.

The two-level or multilevel cache represents a possible approach to the problem of an overlaid single-level cache, but further study is needed.

2.15 Pipelining

Referencing a cache memory is a multistep process. There is the need to obtain priority for a cache cycle. Then the TLB and the desired set are accessed in parallel. After this, the real address is used to select the information in the set, and finally, after the information is read out, the replacement bits are updated. In large, high-speed machines, it is common to pipeline the cache, as well as the rest of the CPU, so that more than one cache access can be in progress at the same time. This pipelining is of various degrees of sophistication, and we illustrate it by discussing two machines: the Amdahl 470V/7 and the IBM 3033.

In the 470V/7, a complete read requires four cycles, known as the P, B1, B2, and R cycles [Smr78b]. The P (priority) cycle is used to determine which of several possible competing sources of requests to the cache will be permitted to use the next cycle. The B1 and B2 (buffer 1, buffer 2) cycles are used actually to access the cache and the TLB, to select the appropriate line from the cache, to check that the contents of the line are valid, and to shift to get the desired byte location out of the two-word (8-byte) segment fetched. The data are available at the end of the B2 cycle. The R cycle is used for "cleanup" and the replacement status is updated at that time. It is possible to have up to four fetches active at any one time in the cache, one in each of the four cycles mentioned above. The time required by a store is longer since it is essentially a read followed by a modify and write-back; it takes six cycles all together, and one store requires two successive cycles in the cache pipeline.

The pipeline in the 3033 cache is similar [IBM79]. The cache in the 3033 can service one fetch or store in each machine cycle, where the turnaround time from initial request for priority until the data is available is about 2½ cycles (½-cycle transmission time to S-unit, 1½ cycles in S-unit, ½ cycle to return data to instruction unit). An important feature of the 3033 is that the cache accesses do not have to be performed in the order that they are issued. In particular, if an access causes a miss, it can be held up while the miss is serviced, and at the same time other requests which are behind it in

the pipeline can proceed. There is an elaborate mechanism built in which prevents this out-of-order operation from producing incorrect results.

2.16 Translation Lookaside Buffer

The translation lookaside buffer (also called the translation buffer [DEC78], the associative memory [Sch71], and the directory lookaside table [IBM78]), is a small, high-speed buffer which maintains the mapping between recently used virtual and real memory addresses (see Figure 2). The TLB performs an essential function since otherwise an address translation would require two additional memory references: one each to the segment and page tables. In most machines, the cache is accessed using real addresses, and so the design and implementation of the TLB is intimately related to the cache memory. Additional information relevant to TLB design and operation may be found in JONE77b, LUDL77, RAMA81, SATY81, SCHN71, and VUK71. Discussions of the use of TLBs (TLB chips or memory management units) in microcomputers can be found in JONH81, SCHN81, and ZOLN81.

The TLB itself is typically designed to look like a small set-associative memory. For example, the 3033 TLB (called the ILAT or directory lookaside table) is set-associative, with 64 sets of two elements each. Similarly, the Amdahl 470V/6 uses 18 sets of two elements each and the 470V/7 and V/8 have 256 sets of 2 elements each. The IBM 3081 TLB has 128 entries.

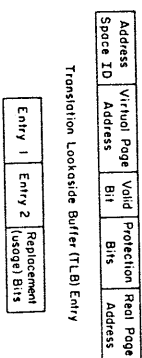
The TLB differs in some ways from the cache in its design. First, for most processes, address spaces start at zero and extend upward as far as necessary. Since the TLB translates page addresses from virtual to real, only the high-order (page number) address bits can be used to access the TLB. If the same method was used as that used for accessing the cache (bit selection using lower order bits), the low-order TLB entries would be used disproportionately and therefore the TLB would be used inefficiently. For this reason, both the 3033 and the 470 hash the address before accessing the TLB (see Figure 2). Consider the 24-bit address used in the System/370, with the address numbered from 1 to 24 (high order to

low order). Then the bits 13 to 24 address the byte within the page (4096-byte page) and the remaining bits (1 to 12) can be used to access the TLB. The 3033 contains a 6-bit index into the TLB computed as follows. Let q be the Exclusive OR operator, a 6-bit quantity is computed $[7, (8 @ 2), (9 @ 3), (10 @ 4), (11 @ 5), (12 @ 6)]$, where each number refers to the input bit it designates.

The Amdahl 470V/7 and 470V/8 use a different hashing algorithm, one which provides much more thorough randomization at the cost of significantly greater complexity. To explain the hashing algorithm, we first explain some other items. The 470 associates with each address space an 8-bit tag field called the address space identifier (ASID) (see Section 2.19.1). We refer to the bits that make up this tag field as S_1, \dots, S_8 . These bits are used in the hashing algorithm as shown below. Also, the 470V/7 uses a different algorithm to hash into each of the two elements of a set; the TLB is more like a pair of direct mapping buffers than a set-associative buffer. The first half is addressed using 8 bits calculated as follows: $[(6 @ 1 @ S_8), 7, (8 @ 3 @ S_6), 9, (10 @ S_4), 11, (12 @ S_2), 5]$; and the second half is addressed as $[6, (7 @ 2 @ S_7), 8, (9 @ 4 @ S_5), 10, (11 @ S_3), 12, (5 @ S_1)]$. There are no published studies that indicate whether the algorithm used in the 3033 is sufficient or whether the extra complexity of the 470V/7 algorithm is warranted.

There are a number of fields in a TLB entry (see Figure 34). The virtual address presented for translation is matched against the virtual address tag field (ASID) plus the virtual page address in the TLB to ensure that the right entry has been found. The virtual address tag field must include the address space identifier (8 bits in the 470V/7, 5 bits in the 3033) so that entries for more than one process can be in the TLB at one time. A protection field (in 370-type machines) is also included in the TLB and is checked to make sure that the access is permissible. (Since keys are associated on a page basis in the 370, this is much more efficient than placing the key with each line in the cache.) The real address corresponding to the virtual address is the primary output of the TLB and occupies a

Figure 34. Structure of translation lookaside buffer (TLB) entry and TLB set.



field. There are also bits that indicate whether a given entry in the TLB is valid and the appropriate bits to permit LRU-like replacement. Sometimes, the modify and reference bits for a page are kept in the TLB. If so, then when the entry is removed from the TLB, the values of those bits must be stored.

It may be necessary to change one or more entries in the TLB whenever the virtual to real address correspondence changes for any page in the address space of any active process. This can be accomplished in two ways: (1) if a single-page table entry is changed (in the 370), the I/PTE (insert page table entry) instruction causes the TLB to be searched, and the now invalid entry purged; (2) if the assignment of address space IDs is changed, then the entire TLB is purged. In the 3033, purging the TLB is slow (16 machine cycles) since each entry is actually invalidated. The 470V/7 does this in a rather clever way. There are two sets of bits used to denote valid and invalid entries, and a flag indicating which set is to be used at any given time. The set not in use is supposed to be set to zero (invalid). The purge TLB command has the effect of flipping this flag, so that the set of bits indicating that all entries are invalid are now in use. The set of bits no longer in use is reset to zero in the background during idle cycles. See Cosc81 for a similar idea.

The cache on the DEC VAX 11/780 [DEC78] is similar to but simpler than that in the IBM and Amdahl machines. A set-associative TLB (called the *translation buffer*) is used, with 64 sets of 2 entries each. (The VAX 11/750 has 256 sets of 2 entries each.) The set is selected by using the high-order address bit and the five low-order bits of the page address, so the address need not be hashed at all. Since the

higher order address bit separates the user from the supervisor address space, this means that the user and supervisor TLB entries never map into the same locations. This is convenient because the user part of the TLB is purged on a task switch. (There are no address space IDs.) The TLB is also used to hold the dirty bit, which indicates if the page has been modified, and the protection key.

Published figures for TLB performance are not generally available. The observed miss ratio for the Amdahl 470V/6 TLB is about 0.3 to 0.4 percent (Private Communication: W. J. Harding). Simulations of the VAX 11/780 TLB [SARV81] show miss ratios of 0.1 to 2 percent for TLB sizes of 64 to 256 entries.

2.17 Translator

When a virtual address must be translated into a real address and the translation does not already exist in the TLB, the *translator* must be used. The translator obtains the base of the segment table from the appropriate place (e.g., control register 1 in 370 machines), adds the segment number from the virtual address to obtain the page table address, then adds the page number (from the virtual address) to the page table address to get the real page address. This real address is passed along to the cache so that the access can be made, and simultaneously, the virtual address/real address pair is entered in the TLB. The translator is basically an adder which knows what to add.

It is important to note that the translator requires access to the segment and page table entries, and these entries may be either in the cache or in main memory. Provision must be made for the translator access to proceed unimpeded, independent of whether target addresses are cache or main memory resident.

We also observe another problem related to translation, "page crossers." The target of a fetch or store may cross from one page to another, in a similar way as for "line crossers." The problem here is considerably more complicated than that of line crossers since, although the virtual addresses are contiguous, the real addresses may not be. Therefore, when a page crossover occurs, two

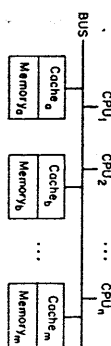


Figure 35. Diagram of computer system in which caches are associated with memories rather than with processors.

separate translations are required; these may occur in the TLB and/or translator as the occasion demands.

2.18 Memory-Based Cache

It was stated at the beginning of this paper that caches are generally associated with the processor and not with the main memory. A different design would be to place the cache in the main memory itself. One way to do this is with a shared bus interfacing between one or more CPUs and several main memory modules, each with its own cache (see Figure 35).

There are two reasons for this approach. First, the access time at the memory module is decreased from the typical 200-500 nanoseconds (given high-density MOS RAMs) to the 50-100 nanoseconds possible for a high-speed cache. Second, there is no consistency problem even though there are several CPUs. All accesses to data in memory module *i* go through cache *i* and thus there is only one copy of a given piece of data.

Unfortunately, the advantages mentioned are not nearly sufficient to compensate for the shortcomings of this design. First, the design is too slow; with the cache on the far side of the memory bus, access time is not cut sufficiently. Second, it is too expensive; there is one cache per memory module. Third, if there are multiple CPUs, there will be memory bus contention. This slows down the system and causes memory access time to be highly variable.

Overall, the scheme of associating the cache with the memory modules is very poor; unless both the main memory and the processors are relatively slow. In that case, a large number of processors could be served by a small number of memory modules with built-in caches, over a fast bus.

2.19 Specialized Caches and Cache Components

This paper has been almost entirely concerned with the general-purpose cache memory found in most large, high-speed computers. There are other caches and buffers that can be used in such machines and we briefly discuss them in this section.

2.19.1 Address Space Identifier Table

In many computers, the operating system identifier for an address space is quite long; in the IBM-compatible machines discussed (370/168, 3033, 470V), the identifier is the contents of control register 1. Therefore, these machines associate a much shorter tag with each address space for use in the TLB and/or the cache. This tag is assigned on a temporary basis by the hardware, and the correspondence between the address space and the tag is held in a hardware table which we name the Address Space Identifier Table (ASIT). It is also called the Segment Base Register Table in the 470V/7, the Segment Table Origin Address Stack in the 3033 [IBM79] and 370/168 [IBM75], and the Segment Base Register Stack in the 470V/6.

The 3033 ASIT has 32 entries, which are assigned starting at 1. When the table becomes full, all entries are purged and IDs are reassigned dynamically as address spaces are activated. (The TLB is also purged.) When a task switch occurs, the ASIT in the 3033 is searched starting at 1; when a match is found with control register 1, the index of that location becomes the address space identifier.

The 470V/6 has a somewhat more complex ASIT. The segment table origin address is hashed to provide an entry into the ASIT. The tag associated with that address is then read out. If the address space does not have a tag, a previously unused tag is assigned and placed in the ASIT. Whenever a new tag is assigned, a previously used tag is made available by deleting its entry in the ASIT and (in the background) purging all relevant entries in the TLB. (A complete TLB purge is not required.) Thirty-two valid tags are available, but the ASIT has the capability of holding up to 128 entries; thus, all 32 valid tags can usually be used, with little fear of hashing conflicts.

2.19.2 Execution Unit Buffers

In some machines, especially the IBM 360/91 [ANDE67b, IBM71, TOMA67], a number of buffers are placed internally in the execution unit to buffer the inputs and outputs of partially completed instructions. We refer the reader to the references just cited for a complete discussion of this.

2.19.3 Instruction Lookahead Buffers

In several machines, especially those without general-purpose caches, a buffer may be dedicated to lookahead buffering of instructions. Just such a scheme is used on the Cray 1 [CHAY76], the CDC 6600 [CDC74], the CDC 7600, and the IBM 360/91 [ANDE67a, BOLAE7, IBM71]. These machines all have substantial buffers, and loops can be executed entirely within these buffers. Machines with general-purpose caches usually do not have much instruction lookahead buffering, although a few extra bytes are frequently fetched. See also BLA280 and KONE80.

2.19.4 Branch Target Buffer

One major impediment to high performance in pipelined computer systems is the existence of branches in the code. When a branch occurs, portions of the pipeline must be flushed and the correct instruction stream fetched. To minimize the effect of these disruptions, it is possible to implement a branch target buffer (BTB) which buffers the addresses of previous branches and their target addresses. The instruction fetch address is matched against the contents of the branch target buffer and if a match occurs, the next instruction fetch takes place from the (previous) target of the branch. The BTB can correctly predict the correct branch behavior more than 90 percent of the time [LRE82]. Something like a branch target buffer is used in the MU-5 [BBE72, MORR79], and the S-1 [MCW77].

2.19.5 Microcode Cache

Many modern computer systems are microcoded and in some cases the amount of microcode is quite large. If the microcode is not stored in sufficiently fast storage, it

is possible to build a special cache to buffer the microcode.

2.19.6 Buffer Invalidation Address Stack (BIAS)

The IBM 370/168 and 3033 both use a store-through mechanism in which any store to main memory causes the line affected to be invalidated in the caches of all processors other than the one which performed the store. Addresses of lines which are to be invalidated are kept in the buffer invalidation address stack (BIAS) in each processor, which is a small hardware implemented queue inside the S-unit. The DEC VAX 11/780 functions in much the same way, although without a BIAS to queue requests. That is, invalidation requests in the VAX have high priority, and only one may be outstanding at a time.

2.19.7 Input/Output Buffers

As noted earlier, input/output streams must be aborted if the processor is not ready to accept or provide data when they are needed. For this reason, most machines have a few words of buffering in the I/O channels or I/O channel controller(s). This is the case in the 370/168 [IBM75] and the 3033 [IBM78].

2.19.8 Write-Through Buffers

In a write-through machine, it is important to buffer the writes so that the CPU does not become blocked waiting for previous writes to complete. In the IBM 3033 [IBM78], four such buffers, each holding a double word, are provided. The VAX 11/780 [DEC78], on the other hand, buffers one write. (Four buffers were recommended in SMIT79.)

2.19.9 Register Cache

It has been suggested that registers be automatically stacked, with the top stack frames maintained in a cache [DIR782]. While this is much better (faster) than implementing registers as part of memory, as with the Texas Instruments 9900 microprocessor, it is unlikely to be as fast as regular, hardware registers. The specific cache described, however, is not general purpose,

but is dedicated to holding registers; it therefore should be much faster than a large, general-purpose cache.

3. D. SECTIONS FOR RESEARCH AND DEVELOPMENT

Cache memories are moderately well understood, but there are problems which indicate directions both for research and development. First, we note that technology is changing; storage is becoming cheaper and faster, as is processor logic. Cost/performance trade-offs and compromises will change with technology and the appropriate solutions to the problems discussed will shift. In addition to this general comment, we see some more specific issues.

3.1 On-Chip Caches and Other Technology Advances

The number of gates that can be placed on a microcomputer chip is growing quickly, and within a few years, it will be feasible to build a general-purpose cache memory on the same chip as the processor. We expect that such an on-chip cache will occur. There is research to be done in designing this within the constraints of the VLSI state of the art. (See Lin781 and Ponn82.)

Cache design is also affected by the implementation technology. MOS VLSI, for example, permits wide associative searches to be implemented easily. This implies that parameters such as set size may change with changing technology. This related aspect of technological change also needs to be studied.

3.2 Multicache Consistency

The problem of multicache consistency was discussed in Section 2.7 and a number of solutions were indicated. Additional commercial implementations are needed, especially of systems with four or more CPUs, before the cost/performance trade-offs can be evaluated.

3.3 Implementation Evaluation

A number of new or different cache designs were discussed earlier, such as the split instruction/data cache, the supervisor/user

cache, the multilevel cache, and the virtual address cache. One or more implementations of such designs are required before their desirability can be fully evaluated.

3.4 Hit Ratio versus Size

There is no generally accepted model for the hit ratio of a cache as a function of its size. Such a model is needed, and it will probably have to be made specific to each machine architecture and workload type (e.g., 370 commercial, 370 scientific, and PDP-11).

3.5 TLB Design

A number of different TLB designs exist, but there are almost no published evaluations (but see SAT781). It would be useful to know what level of performance can be expected from the various designs and, in particular, to know whether the complexity of the Amdahl TLB is warranted.

3.6 Cache Parameters versus Architecture and Workload

Most of the studies in this paper have been based on IBM System 370 user program address traces. On the basis of that data, we have been able to suggest desirable parameter values for various aspects of the cache. Similar studies need to be performed for other machines and workloads.

APPENDIX. EXPLANATION OF TRACE NAMES

1. EDC PDP-11 trace of text editor, written in C, compiled with C compiler on PDP-11.
2. ROFPAS PDP-11 trace of text output and formatting program. (called ROFP or runoff).
3. TRACE PDP-11 trace of program tracer itself tracing EDC. (Tracer is written in assembly language.)
4. FGO1 FORTRAN Go step, factor analysis (1249 lines, single precision).
5. FGO2 FORTRAN Go step, double-precision analysis of satellite information, 2057 lines, FortG compiler.
6. FGO3 FORTRAN Go step, double-

- precision numerical analysis, 840 lines, FortG compiler.
7. FGO4 FORTRAN Go step, FFT of hole in rotating body. Double-precision FortG.
8. CGO1 COBOL Go step, fixed-assets program doing tax transaction selection.
9. CGO2 COBOL Go step, "fixed assets: year end tax select."
10. CGO3 COBOL Go step, projects depreciation of fixed assets.
11. PGO2 PL/I Go step, does CCW analysis.
12. IEBDG IBM utility that generates test data that can be used in program debugging. It will create multiple data sets of whatever form and contents are desired.
13. PGO1 PL/I Go step, SMF billing program.
14. FCOMF FORTRAN compile of program that solves Reynolds partial differential equation (2330 lines).
15. CCOMP COBOL compile, 240 lines, accounting report.
16. WATEX Execution of a FORTRAN program compiled using the WATFIV compiler. The program is a combinatorial search routine.
17. WATFIV FORTRAN compilation using the WATFIV compiler. (Compiles the program whose trace is the WATEX trace.)
18. APL Execution of APL program which does plots at a terminal.
19. FFT Execution of an FFT program written in ALGOL, compiled using ALGOLW compiler at Stanford.

ACKNOWLEDGMENTS

This research has been partially supported by the National Science Foundation under grants MCS77-28429 and MCS-8202591, and by the Department of Energy under contract FY-76-C-03-0615 (with the Stanford Linear Accelerator Center).

I would like to thank a number of people who were of help to me in the preparation of this paper. Matt Diehl, George Rossman, Alan Knowles, and Dileep Bhandarkar helped me to obtain materials describing various machines. George Rossman, Bill Harding, Jim Goodman, Domenico Ferrai, Mark Hill, and Bob

Doran read and commented upon a draft of this paper. Mac MacDougall provided a number of other useful comments. John Lee generated a number of the traces used for the experiments in this paper from trace data available at Andahl Corporation. Four of the traces were generated by Len Shumick. The responsibility for the contents of this paper, of course, remains with the author.

REFERENCES

- BARSAIMAN, H., AND DECCECAM, A. "System design considerations of cache memories." in *Proc. IEEE Computer Society Conference* (1972), IEEE, New York, pp. 107-110.
- BEAN, B. M., LANGSTON, K., PART-
RIDGE, R. S., K. B. Bias filter memory
for filtering out unnecessary interroga-
tions of cache directories in a multi-
processor system. United States Patent 4-
142,224, Feb. 17, 1978.
- BENNETT, S. Cache management
system using virtual and real tags in the
cache directory. *IBM Tech. Disclosure
Bull.* 21, 11 (April 1978), 4641.
- BRAY, L. A. A study of replacement
algorithms for a virtual storage com-
puter. *IBM Syst. J.* 6, 2 (1966), 78-101.
- BRILL, J., CAMASANT, D., AND BRILL, G.
An investigation of alternative
cache organizations. *IEEE Trans. Com-
put. TC-23*, 4 (April 1974), 346-351.
- BRINKERT, B. T., AND FRANKCZAK, P.
A. Cache memory with prefetching of
data by priority. *IBM Tech. Disclosure
Bull.* 18, 12 (May 1976), 4231-4232.
- BENNETT, B. T., POWERS, J. H., P-
ZAK, T. R., AND RECHTSCHAFFEN, R.
N. Prefetching in a multilevel memory
hierarchy. *IBM Tech. Disclosure Bull.*
25, 1 (June 1982), 88.
- BENG, H. S., AND SUMMERFIELD, A.
R. CPU busy not all productive uti-
lization. *Share Computer Measurement
and Evaluation Newsletter*, no. 39
(Sept. 1976), 95-97.
- BENGEY, A. L., JR. Increased computer
throughput by conditioned memory data
prefetching. *IBM Tech. Disclosure Bull.*
20, 10 (March 1978), 4103.
- BLAZEWSKI, T. J., DOBRZYNSKI, S. M.,
AND WATSON, W. D. Instruction buffer
with simultaneous storage and fetch op-
erations. *IBM Tech. Disclosure Bull.* 23,
2 (July 1980), 670-672.
- BLOUNT, F. T., BURTON, R. J., MAR-
TIN, D. B., MCGILVER, B. L., AND ION-
BINSON, J. R. Deferred cache storing
method. *IBM Tech. Disclosure Bull.* 23,
1 (June 1980), 262-263.
- BOLAND, L. J., GRANT, G. D., MAR-
COTTE, A. V., MESSINA, B. U., AND
SMITH, J. W. The IBM System/360
Model 91: Storage system. *IBM J. Res.
Dev.* 11, 1 (Jan 1967), 54-68.
- BORGERSON, B. R., GODFREY, M. D.,
HAGERTY, P. E., RYKSEN, T. R. "The
architecture of the Sperry Univac 1100
series systems." in *Proc. 6th Annual
Symp. Computer Architecture* (April 23-
25, 1979), ACM, New York, N.Y., pp.
137-146.
- CAMPBELL, J. E., STROHM, W. G., AND
TEMPLE, J. L. Most recent class used
search algorithm for a memory cache.
IBM Tech. Disclosure Bull. 18, 10
(March 1976), 3307-3308.
- CONTROL DATA CORP. Control Data
6800 Series Computer Systems Refer-
ence Manual, Arden Hills, Minn., 1974.
- CRAIG, L., AND FEATURNER, P. A
new solution to coherence problems in
multicache systems. *IEEE Trans. Com-
put. TC-22*, 12 (Dec. 1978), 1112-1118.
- CHU, D. K. Optimum implementation
of LRU hardware for a 4-way set-asso-
ciative memory. *IBM Tech. Disclosure
Bull.* 17, 11 (April 1975), 3161-3163.
- CHOW, C. K. Determining the optimum
capacity of a cache memory. *IBM Tech.
Disclosure Bull.* 17, 10 (March 1975),
3163-3168.
- CHOW, C. K. Determination of cache's
capacity and its matching storage hier-
archy. *IEEE Trans. Comput. TC-25*, 2
(Feb. 1976), 157-164.
- CHU, W. W., AND OPPERHECK, H.
Program behavior and the page
fault frequency replacement algorithm.
Computer 9, 11 (Nov. 1976), 29-38.
- CLARK, D. W., LAMSON, B. W., PIER,
K. A. The memory system of a high
performance personal computer. *IEEE
Trans. Comput. TC-30*, 10 (Oct. 1981),
715-733.
- CLARK, D. W. Cache performance in
the VAX-11/780. To appear in *ACM
Trans. Comp. Syst.* 1, 1 (Feb. 1983).
- CORFMAN, E. G., AND DENNING, P. J.
Operating Systems Theory. Pren-
tice-Hall, Englewood Cliffs, N.J., 1973.
- CONRI, C. J., GIBSON, D. H., AND PIR-
KOWSKY, S. H. Structural aspects of
the system/360 Model 85. *IBM Syst. J.*
7, 1 (1969), 2-21.
- CONRI, C. J. Concepts for buffer stor-
age. *IEEE Computer Group News* 2, 8
(March 1969), 9-13.
- COSCARIELLA, A. S., AND SKILLERS, F.
F. System for purging TLB. *IBM Tech.
Disclosure Bull.* 24, 2 (July 1981), 910-
911.
- CRAV RESEARCH, INC. Cray-1 Com-
puter System Reference Manual, Bloom-
ington, Minn., 1976.
- DEC CORP. Equipment Corp. "TJB/
Digital Control Technical Descrip-
tion—Vax-11/780 Implementation." Document No. EK-NM780-TJD-001,
First Edition (April 1978), Digital Equip-
ment Corp., Maynard, Mass., 1978.
- DENNING, P. J. The working set model
for program behavior. *Commun. ACM*
11, 5 (May 1968), 323-333.
- DENNING, P. J. "On modeling program
behavior." in *Proc. Spring Joint Com-
puter Conference*, vol. 40, AFIPS Press,
Arlington, Va., 1972, pp. 937-944.
- DIERHELM, M. A. "Level 66 cache
memory." Tech. Info. Notepad I-114,
Honeywell, Phoenix, Ariz., April, 1974.
- DITZEN, D. R. "Register allocation for
free: The C machine stack cache." in
*Proc. Symp. on Architectural Support
for Programming Languages and Op-
erating Systems* (Tulane Univ., Calif.,
March 1-3, 1982), ACM, New York,
N.Y., 1982.
- DRIMAK, E. G., DUTTON, P. F., HICKS,
G. L., AND STRIER, W. R. Multi-
processor locking with a bypass for chan-
nel references. *IBM Tech. Disclosure
Bull.* 23, 12 (May 1981), 5329-5331.
- DRIMAK, E. G., DUTTON, P. F., AND
STRIER, W. R. Attached processor si-
multaneous data searching and transfer
via main storage controls and intercache
transfer controls. *IBM Tech. Disclosure
Bull.* 24, 1A (June 1981), 26-27.
- DRISCOLL, G. C., MATTHE, R. E., PIZAK,
T. R., AND SHEDLETSKY, J. J. Split
cache with variable intercache boundary.
IBM Tech. Disclosure Bull. 22, 11 (April
1980), 5183-5186.
- DUNOIS, M., AND BRIGGS, F. A.
"Effects of cache concurrency in multi-
processors." in *Proc. 9th Annual Symp.
Computer Architecture* (Austin, Texas,
April, 1982), ACM, New York, N.Y.,
1982, pp. 292-308.
- EASTON, M. C., AND FAGIN, R. "Cold-
start vs. warm-start miss ratios and mul-
tiprogramming performance." *IBM Res.
Rep. RC 5175*, Nov., 1976.
- EASTON, M. C. Computation of cold
start miss ratios. *IEEE Trans. Comput.*
TC-27, 5 (May 1978), 404-408.
- ELECTRONICS MAGAZINE. Altering
computer architecture is way to raise
throughput, suggests IBM researchers.
Dec. 23, 1976, 30-31.
- ELECTRONICS New TI 16-bit machine
has on-chip memory. Nov. 3, 1981, 57.
- ENGEL, T. A. Paged control store pre-
fetch mechanism. *IBM Tech. Disclosure
Bull.* 16, 7 (Dec. 1973), 2140-2141.
- FAYNE, P., AND KUHN, R. Fast mem-
ory organization. *IBM Tech. Disclosure
Bull.* 21, 2 (July 1978), 649-650.
- FORUKAWA, K., AND KASAI, T. "The
efficient use of buffer storage." in *Proc.
ACM 1977 Annual Conference* (Seattle,
Wa., Oct. 16-19, 1977), ACM, New York,
N.Y., pp. 399-403.
- FORNEY, R. W. Selection of least re-
cently used slot with bad entry and
locked slots involved. *IBM Tech. Dis-
closure Bull.* 21, 6 (Nov. 1978), 230.
- GAGNER, J. Determining hit ratios for
multilevel hierarchies. *IBM J. Res. Dev.*
18, 4 (July 1974), 316-327.
- GIBSON, D. H. "Consideration in block-
oriented systems design." in *Proc.
Spring Joint Computer Conf.*, vol. 30,

- Thompson Books, Washington, D.C., 1967, pp. 75-80.
- GIND77 Gindrel, J. D. Buffer block prefetching method. *IBM Tech. Disclosure Bull.* 20, 2 (July 1977), 696-697.
- GREB74 Greenberg, B. S. "An experimental analysis of program reference patterns in the multivirt memory." Project MAC Tech. Rep. MAC-TR-127, 1974.
- GUST82 Gustafson, R. N., and Spafacio, F. J. IBM 3081 processor unit: Design considerations and design process. *IBM J. Res. Dev.* 26, 1 (Jan. 1982), 12-21.
- HAL79 Halpern, B., and Hirsom, B. "S-1 architecture manual." Tech. Rep. No. 161, Computer Systems Laboratory, Stanford Univ., Stanford, Calif., Jan., 1979.
- HARD75 Harding, W. J. "Hardware Controlled Memory Hierarchies and Their Performance." Ph.D. dissertation, Arizona State Univ., Dec., 1975.
- HARD80 Harding, W. J., MacDougall, M. H., Raymond, W. J. "Empirical estimation of cache miss ratios as a function of cache size." Tech. Rep. PN 820-420-700A (Sept. 26, 1980), Amдах Corp.
- HORE81a Horev, L. W., and Voldman, J. Mechanism for cache replacement and prefetching driven by heuristic estimation of operating system behavior. *IBM Tech. Disclosure Bull.* 23, 8 (Jan. 1981), 3923.
- HORE81b Horev, L. W., and Voldman, J. Cache line reclamation and cast out avoidance under operating system control. *IBM Tech. Disclosure Bull.* 23, 8 (Jan. 1981), 3912.
- LABE72 Labert, R. The MVS instruction pipeline. *Comput. J.* 16, 1 (Jan. 1972), 42-50.
- LABE77 Labert, R. N., and Husand, M. A. The MVS name store. *Comput. J.* 20, 3 (Aug. 1977), 227-231.
- IBM71 IBM "IBM system/360 and System/370 Model 165 Functional characteristics." Form GA22-6945-2 (Nov. 1971), IBM, Armonk, N.Y.
- IBM75 IBM "IBM System/370 Model 168 Theory of Operation/Diagrams Manual—Processor Storage Control Function (PSCF)." vol. 4, IBM, Poughkeepsie, N.Y., 1975.
- IBM78 IBM "3033 Processor Complex, Theory of Operation/Diagrams Manual—Processor Storage Control Function (PSCF)." vol. 4, IBM, Poughkeepsie, N.Y., 1978.
- IBM82 IBM "Form GA22-7076, IBM, Poughkeepsie, N.Y., 1982.
- JOHN81 Johnson, R. C. Microsystems exploit mainframe methods. *Electronics*, Aug. 11, 1981, 119-127.
- JONE76 Jones, J. D., Junon, D. M., Parringer, R. L., and Shawley, B. L. Updating cache data arrays with data stored by other CPUs. *IBM Tech. Disclosure Bull.* 19, 2 (July 1976), 594-596.
- JONE77a Jones, J. D., and Junon, D. M. Cache address directory invalidation scheme for multiprocessing system. *IBM Tech. Disclosure Bull.* 20, 1 (June 1977), 295-296.
- JONE77b Jones, J. D., and Junon, D. M. Prefetch lookahead buffer. *IBM Tech. Disclosure Bull.* 20, 1 (June 1977), 297-298.
- KAP73 Kaplan, K. R., and Windsor, R. O. Cache-based computer systems. *IEEE Computer* 6, 3 (March 1973), 40-56.
- KOBA73 Kobayashi, M. "An algorithm to measure the buffer growth function." Tech. Rep. PN 820413-700A (Aug. 8, 1980), Amдах Corp.
- KON80 Konzen, D. H., Martin, D. B., McGivray, B. L., and Tokasuro, H. M. Demand driven instruction fetching inhibit mechanism. *IBM Tech. Disclosure Bull.* 23, 2 (July 1980), 716-717.
- KROFT, D. "Lockup-free instruction fetch/prefetch cache organization," in *Proc. 8th Annual Symp. Computer Architecture* (Minneapolis, Minn., May 12-14, 1981), ACM, New York, N.Y., pp. 81-87.
- KUMAR, B. "A model of spatial locality and its application to cache design." Tech. Rep. (unpubl.), Computer Systems Laboratory, Stanford Univ., Stanford, Calif., 1979.
- LAFRE77 LaFrance, D. S., and Guttag, K. M. Fast on-chip memory extends 16 bit finline's reach. *Electronics*, Feb. 24, 1981, 157-161.
- LAWSON, B. W., and Pien, K. A. "A processor for a high-performance personal computer," in *Proc. 7th Annual Symp. Computer Architecture* (May 6-8, 1980), ACM, New York, N.Y., pp. 146-150.
- LEE, F. F. Study of "look-aside" memory. *IEEE Trans. Comput.* TC-18, 11 (Nov. 1969), 1062-1064.
- LEE, J. M., and Weinberger, A. A solution to the synonym problem. *IBM Tech. Disclosure Bull.* 22, 8A (Jan. 1980), 3331-3333.
- LEE, J. K. F., and Smith, A. J. "Analysis of branch prediction strategies and branch target buffer design." Tech. Rep., Univ. of Calif., Berkeley, Calif., 1982.
- LEHMAN, A., and Schmid, D. "The performance of small cache memories in minicomputer systems with several processors," in *Digital Memory and Storage*, Springer-Verlag, New York, 1978, pp. 391-407.
- LEHMAN, A. Performance evaluation and prediction of storage hierarchies. Source unknown, 1980, pp. 43-54.
- LEWIS, P. A. W., and Yue, P. C. "Statistical analysis of program reference patterns in a paging environment," in *Proc. IEEE Computer Society Conference*, IEEE, New York, N.Y., 1971.
- LEWIS, P. A. W., and Shredler, G. S. Empirically derived micro models for sequences of page exceptions. *IBM J. Res. Dev.* 17, 2 (March 1973), 80-100.
- LINDSAY, D. C. Cache memories for microprocessors. *Computer Architecture News* 9, 5 (Aug. 1981), 6-13.
- LITRAY, J. S. Structural aspects of the System/360 Model 65, II the cache. *IBM Syst. J.* 7, 1 (1968), 15-21.
- LIU, L. Cache-splitting with information of XT-sensitivity in tightly coupled multiprocessing systems. *IBM Tech. Disclosure Bull.* 25, 1 (June 1982), 54-55.
- LOSG, J. J., Parks, L. S., Sachar, H. E., and Yamoun, J. Conditional cache miss facility for handling short/long cache requests. *IBM Tech. Disclosure Bull.* 25, 1 (June 1982), 110-111.
- LUTON, D. M., and Moore, B. B. Channel DAT with pin bus. *IBM Tech. Disclosure Bull.* 20, 2 (July 1977), 683.
- MACDOUGALL, M. H. "The stack growth function model." Tech. Rep. 820228-700A (April 1979), Amдах Corp.
- MARUYAMA, K. mLRU page replacement algorithm in terms of the reference matrix. *IBM Tech. Disclosure Bull.* 17, 10 (March 1975), 3101-3103.
- MARUYAMA, K. Implementation of the stack operation circuit for the LRU algorithm. *IBM Tech. Disclosure Bull.* 19, 1 (June 1976), 321-323.
- MATTHEW, R. L. Evaluation of multilevel memories. *IEEE Trans. Magnetics* MAG-7, 4 (Dec. 1971), 814-819.
- MATTHEW, R. L., Gershy, J., Slutz, D. R., and Traicors, I. L. Evaluation techniques for storage hierarchies. *IBM Syst. J.* 9, 2 (1970), 78-117.
- MAZARE, G. "A few examples of how to use a symmetrical multi-micro-processor," in *Proc. 4th Annual Symp. Computer Architecture* (March 1977), ACM, New York, N.Y., pp. 57-62.
- MCWILLIAMS, T., Winrows, L. C., and Wood, L. "Advanced digital processor technology base development for many applications: The S-1 Processor." Tech. Rep. UCIO-117705, Lawrence Livermore Laboratory, Sept. 1977.
- MEAD, R. M. "On memory system design," in *Proc. Fall Joint Computer Conference*, vol. 37, AFIPS Press, Arlington, Va., 1970, pp. 33-43.
- MILANING, G., and Mirkov, R. "VSI-12 experience at the University of Toronto Computer Centre," in *Share 44 Proc.* (Los Angeles, Calif., March, 1975), pp. 1887-1895.
- MORRIS, D., and Labert, R. N. The MVS Computer System. Springer-Verlag, New York, 1979.
- NEAL, C. H., and Wassser, E. R. Shadow directory for attached processor system. *IBM Tech. Disclosure Bull.* 23, 8 (Jan. 1981), 3667-3668.
- NEAL, C. H., and Wassser, E. R. Two-level DIAT hierarchy. *IBM Tech. Disclosure Bull.* 24, 9 (Feb. 1982), 4714-4715.
- OHNO, N., and Hakozi, K. Pseudo random access memory system with CCD-SRAM and MOS RAM on a chip. 1977.
- OLBERT, A. G. Fast DIAT load for V-R translations. *IBM Tech. Disclosure Bull.* 22, 4 (Sept. 1979), 1434.
- PERKINS, D. R. "The Design and Management of Predictive Caches." Ph.D. dissertation, Univ. of Calif., San Diego, Calif., 1980.
- PERUTTO, B. L., and Shustek, L. J. "An instruction timing model of CPU performance," in *Proc. 4th Annual Symp. Computer Architecture* (March 1977), ACM, New York, N.Y., pp. 163-178.
- POIM, A. V., Acharya, O. P., Cheng, C.-W., and Shum, A. C. An efficient flexible buffered memory system. *IEEE Trans. Magnetics* MAG-9, 3 (Sept. 1973), 173-179.
- POIM, A. V., Acharya, O. P., and Monroze, R. N. "The cost and performance tradeoffs of buffered memories," in *Proc. IEEE* 63, 8 (Aug. 1975), pp. 1129-1135.
- POIM, A. V., and Acharya, O. P. "A cache technique for bus oriented multiprocessor systems," in *Proc. Computer* (San Francisco, Calif., Feb. 1982), 1535E, New York, pp. 62-65.
- POIM, J. H., and Rechtshaffen, R. N. Base/displacement lookahead buffer. *IBM Tech. Disclosure Bull.* 22, 11 (April 1980), 5182.
- POWELL, M. L. "The DEMOS File system," in *Proc. 6th Symp. on Operating Systems Principles* (West Lafayette, Ind., Nov. 16-18, 1977), ACM, New York, N.Y., pp. 33-42.
- RABIN, G. M. "The 801 minicomputer," in *Proc. Symp. on Architectural Support for Programming Languages and Operating Systems* (Palo Alto, Calif., March 1-3, 1982), ACM, New York, N.Y., pp. 39-47.
- RAMMOHANARAO, K., and Sacks-Davies, R. Hardware address translation for machines with a large virtual memory. *Inf. Process. Lett.* 13, 1 (Oct. 1981), 23-29.
- RAU, B. R. "Sequential prefetch strategies for instructions and data." Digital

- systems laboratory tech. rep. 131 (1976), Stanford Univ., Stanford, Calif.
- REILLY, J., SUTTON, A., NASSER, R., AND GILSON, R. Processor controller for the IBM 3081. *IBM J. Res. Dev.* 26, 1 (Jan. 1982), 22-29.
- RISE, F. N., AND WARREN, H. S., JR. Read-constant control line to cache. *IBM Tech. Disclosure Bull.* 20, 6 (Nov. 1977), 2509-2510.
- ROSEMAN, G. Private communication. Palyn Associates, San Jose, Calif., 1979.
- SALTZ, J. H. "A simple linear model of demand paging performance. *Commun. ACM* 17, 4 (April 1974), 181-186.
- SAVANAVAYAN, M., AND BHANDARKAR, D. Design trade-offs in VAX-11 translation buffer organization. *IEEE Computer* (Dec. 1981), 103-111.
- SCHROEDER, M. D. "Performance of the GE-645 associative memory while multibits is in operation." in *Proc. 1971 Conference on Computer Performance Evaluation* (Harvard Univ., Cambridge, Mass.), pp. 227-245.
- SHEDDEN, G. S., AND SLUTZ, D. R. Derivation of rates ratios for merged access streams. *IBM J. Res. Dev.* 20, 5 (Sept. 1976), 505-517.
- SLUTZ, D. R., AND TRAIGER, I. L. "Evaluation techniques for cache memory hierarchies." IBM Res. Rep. RJ 1045, May 1972.
- SMITH, A. J. A comparative study of set associative memory mapping algorithms and their use for cache and main memory. *IEEE Trans. Softw. Eng.* SE-4, 2 (March 1978), 121-130.
- SMITH, A. J. Sequential program prefetching in memory hierarchies. *IEEE Computer* 11, 12 (Dec. 1978), 7-21.
- SMITH, A. J. Sequentiality and prefetching in database systems. *ACM Trans. Database Syst.* 3, 3 (Sept. 1978), 223-247.
- SMITH, A. J. Bibliography on paging and related topics. *Operating Systems Review* 12, 4 (Oct. 1978), 39-56.
- SMITH, A. J. Characterizing the storage process and its effect on the update of main memory by write-through. *J. ACM* 26, 1 (Jan. 1979), 6-27.
- SNOW, E. A., AND SIEMOREK, D. P. "Impact of implementation design tradeoffs on performance: The PDP-11. A case study." Dep. of Computer Science Report (Feb. 1978). Carnegie-Mellon University, Pittsburgh, Pa.
- SPARACIO, F. J. Data processing system with second level cache. *IBM Tech. Disclosure Bull.* 21, 6 (Nov. 1978), 2468-2469.
- STEVENSON, D. "Virtual memory on the Z8003." in *Proc. IEEE Comput.* (San Francisco, Calif., Feb. 1981), pp. 355-357.
- STRUCKER, W. D. "Cache memories for PDP-11 family computers." in *Proc. 3rd Annual Symp. Computer Architecture* (Jan. 19-21, 1976). ACM, New York, N.Y., pp. 155-158.
- TANG, C. K. "Cache system design in the tightly coupled multiprocessor system." in *Proc. AFIPS National Computer Conference* (New York City, New York, June 7-10, 1976), vol. 45, AFIPS Press, Arlington, Va., pp. 749-753.
- THAKKAR, S. S. "Investigation of Buffer Store Organization." Master's of science thesis, Victoria University of Manchester, England, October, 1978.
- TOMASULO, R. M. An efficient algorithm for exploiting multiple arithmetic units. *IBM J. Res. Dev.* 11, 1 (Jan. 1967), 25-33.
- WILKES, M. V. Store memories and segmentation. *IEEE Trans. Comput.* (June 1971), 674-675.
- WINDER, R. O. A data base for computer performance evaluation. *IEEE Computer* 6, 3 (March 1973), 25-29.
- YAMOUR, J. Odd/even interleave cache with optimal hardware array cost, cycle time and variable data part width. *IBM Tech. Disclosure Bull.* 23, 7B (Dec. 1980), 3461-3463.
- YEN, W. C., AND FU, K. S. "Analysis of multiprocessor cache organizations with alternative main memory update policies." in *Proc. 8th Annual Symp. Computer Architecture* (Minneapolis, Minn., May 12-14, 1981). ACM, New York, N.Y., pp. 89-105.
- YUVA, A. "165/H58 analysis." Share Inc., Computer Measurement and Evaluation, Selected Papers from the Share Project, vol. III, pp. 585-606, 1975.
- ZOLNOWSKY, J. "Philosophy of the M068451 memory management unit." in *Proc. IEEE Comput.* (San Francisco, Calif., Feb. 1981). IEEE, New York, pp. 358-361.

BIBLIOGRAPHY

- ACKLAND, B. D., AND PUCKENELL, D. A. Studies of cache store behavior in a real time multiprocessor environment. *Electronics Letters* 11, 24 (Nov. 1975), 588-590.
- ACKLAND, B. D. "A bit-slice cache controller." in *Proc. 6th Annual Symp. Computer Architecture* (April 23-25, 1979). ACM, New York, N.Y., pp. 75-82.
- AGRAWAL, O. P., ZINCO, R. J., PONT, A. V. "Applicability of cache memories to dedicated multiprocessor systems." in *Proc. IEEE Computer Society Confer-*

- ence (San Francisco, Calif., Spring 1977). IEEE, New York, pp. 74-76.
- AMDahl, Corp. 470V/6 Machine Reference Manual, 1976.
- BEILL, C. G., AND CASASSENT, D. Implementation of a buffer memory in minicomputers. *Comput. Des.* 10, (Nov. 1971), 83-88.
- BLOOM, L., COHEN, M., AND PORTER, S. "Considerations in the design of a computer with high logic-to-memory speed ratio." in *Proc. Gigacycle Computing Systems* (Jan. 1962). AIEE Special Publication S-136, pp. 53-63.
- BOROWANDY, P. A. "Cache structures based on the execution stack for high level languages." Tech. Rep. 81-08-04, Dep. of Computer Science, Univ. of Washington, Seattle, Wa., 1981.
- BURICK, F. A., AND DUNN, M. "Performance of cache-based multiprocessors." in *Proc. ACM/SIGMETRICS Conf. on Measurement and Modeling of Computer Systems* (Las Vegas, Nev., Sept. 14-16, 1981). ACM, New York, N.Y., 181-190.
- CANNON, J. W., GRIMES, D. W., AND HERMANN, B. D. Storage protect operations. *IBM Tech. Disclosure Bull.* 24, 2 (July 1981), 1184-1186.
- CAPORSA, R. S., DEVERA, J. A., HELLER, A. R., AND MESSITT, J. W. Dynamic address translator for I/O channels. *IBM Tech. Disclosure Bull.* 23, 12 (May 1981), 5503-5508.
- CASPER, P. G., FAIX, M., GOETZ, V., AND ULLAND, H. Cache resident processor registers. *IBM Tech. Disclosure Bull.* 22, 6 (Nov. 1979), 2317-2318.
- CONDROO, H., AND CHAMBERS, J. B. Second group of IBM 4341 machines outdoes the first. *Electronics* (April 7, 1981), 149-152.
- ELLER, J. E., III. "Cache design and the X-tree high speed memory buffers." Master's of science project report, Computer Science Division, EECS Department, Univ. of Calif., Berkeley, Calif., Sept., 1979.
- FARIS, S. M., HENKEL, W. H., VAISAMAKIS, E. A., AND ZAPPE, H. H. Basic design of a Josephson technology cache memory. *IBM J. Res. Dev.* 24, 2 (March 1980), 143-154.
- FARMER, D. Comparing the 4341 and M80/42. *Computerworld*, Feb. 9, 1981.
- FERRICHI, R. A., AND SACHAR, H. E. Interleaved multiple speed memory controls with high speed buffer. *IBM Tech. Disc. Bull.* 22, 5 (Octo. 1979), 1999-2000.
- GARCIA, L. C. Instruction buffer design. *IBM Tech. Disclosure Bull.* 20, 11b (April 1978), 4832-4833.
- HARTEN, R. L., AND MOYER, J. T. Split cycle for a shared data buffer array. *IBM Tech. Disclosure Bull.* 21, 6 (Nov. 1978), 2293-2294.
- HOFFMAN, R. L., MITCHELL, G. R., AND SOUTHS, F. G. Reference and change bit recording. *IBM Tech. Disclosure Bull.* 23, 12 (May 1981), 5516-5519.
- HURSTICH, J., AND SUTLER, W. R. Cache reconfiguration. *IBM Tech. Disclosure Bull.* 23, 9 (Feb. 1981), 4117-4118.
- IBM. "IBM field engineering theory of operation, System/360, Model 195, storage control unit buffer storage," first edition (Aug. 1970). IBM, Poughkeepsie, N.Y.
- IZUKA, H., AND TERU, T. Cache memory simulation by a new method of address pattern generation. *J. IRE* 14, 9 (1972), 663-676.
- KNEPPER, R. W. Cache bit selection circuit. *IBM Tech. Disclosure Bull.* 22, 1 (June 1979), 142-143.
- KNOKE, P. "An analysis of buffered memories." in *Proc. 2nd Hawaii Int. Conf. on System Sciences* (Jan. 1969), pp. 397-400.
- KOROK, A. "Lecture notes for CS252." of course notes (Spring 1976). Univ. of Calif., Berkeley, Calif., 1976.
- LANCER, R. E., DIETRICH, M. A., ISHMAVI, P. C. Cache memory store in a processor of a data processing system. United States Patent 3,986,419, July, 1976.
- LARKER, R. A., LASSETTER, E. R., MOORE, B. B., AND STRICKLAND, J. P. Channel DAT and page pinning for block unit transfers. *IBM Tech. Disclosure Bull.* 23, 2 (July 1980), 704-705.
- LEE, P. A., GHANI, N., AND HERON, K. A recovery cache for the PDP-11. *IEEE Trans. Comput.* TC-29, 6 (June 1980), 546-549.
- LORIN, H., AND GOLDSTEIN, B. "An inversion of the memory hierarchy." IBM Res. Rep. RC 8171, March, 1980.
- MADDOCK, R. F., MARKS, B. L., MINSHULL, J. F., AND PRINCE, M. C. Hardware address relocation for variable length segments. *IBM Tech. Disclosure Bull.* 23, 11 (April 1981), 5186-5187.
- MARTIN, J. R., MAYFIELD, M. J., AND ROWLAND, H. E. Reference associative cache mapping. *IBM Tech. Disclosure Bull.* 23, 9 (Feb. 1981), 3683-3691.
- MEAR, R. M. Design approaches for cache memory control. *Comput. Des.* 10, 1 (Jan. 1971), 87-93.
- MERRILL, B. 370/168 cache memory performance. *Share Computer Measurement and Evaluation Newsletter*, no. 25 (July 1974), 98-101.

- Moore, B. B., RodriL, J. T., Sutton, A. J., and Vowell, J. D. Vary storage physical on/off line in a non-store-through cache system. *IBM Tech. Disclosure Bull.* 23, 7B (Dec. 1980), 3329.
- Nakamura, T., Hagihara, H., Kirigawa, H., and Kanazawa, M. Simulation of a computer system with buffer memory. *J. IPSJ* 15, 1 (1974), 26-33.
- Ngai, C. H., and Strien, W. R. Two-bit DLAT LRU algorithm. *IBM Tech. Disclosure Bull.* 22, 10 (March 1980), 4488-4490.
- Rao, G. S. "Performance analysis of cache memories." Digital Systems Laboratory Tech. Rep. 110 (Aug. 1975), Stanford Univ., Stanford, Calif.
- Razitschaffren, R. N. Using a branch history table to prefetch cache lines. *IBM Tech. Disclosure Bull.* 22, 12 (May 1980), 5539.
- Schulz, C. Fast address translation in systems using virtual addresses and a cache memory. *IBM Tech. Disclosure Bull.* 21, 2 (Jan. 1978), 663-664.
- Thiruvirt, B. A VLSI approach to cache memory. *Comput. Res.* (Jan. 1982), 169-173.
- Trendenick, H. L., and Welch, T. A. "High-speed buffering for variable length operands," in *Proc. 4th Annual Symp. Computer Architecture* (March 1977), ACM, New York, N.Y., pp. 205-210.
- Voldman, J., and Hoevel, L. W. "The fourier cache connection," in *Proc. IEEE Comput. (San Francisco, Calif., Feb. 1981)*, IEEE, New York, pp. 344-354.
- Voldman, J., and Hoevel, L. W. The software-cache connection. *IBM J. Res. Dev.* 25, 6 (Nov. 1981), 877-893.
- Whites, M. V. Slave memories and dynamic storage allocation. *IEEE Trans. Comput.* TC-14, 2 (April 1965), 270-271.

Received December 1979; final revision accepted January 1982.