

Page Placement Algorithms for Large Real-Indexed Caches

R. E. KESSLER and MARK D. HILL

University of Wisconsin

When a computer system supports both paged virtual memory and large real-indexed caches, cache performance depends in part on the main memory page placement. To date, most operating systems place pages by selecting an arbitrary page frame from a pool of page frames that have been made available by the page replacement algorithm. We give a simple model that shows that this naive (arbitrary) page placement leads to up to 30% unnecessary cache conflicts. We develop several page placement algorithms, called *careful-mapping algorithms*, that try to select a page frame (from the pool of available page frames) that is likely to reduce cache contention. Using trace-driven simulation, we find that careful mapping results in 10–20% fewer (dynamic) cache misses than naive mapping (for a direct-mapped real-indexed multimegabyte cache). Thus, our results suggest that careful mapping by the operating system can get about half the cache miss reduction that a cache size (or associativity) doubling can.

Categories and Subject Descriptors: B.3.2 [Memory Structures]: Design Styles—*cache memories, virtual memory*; B.3.3 [Memory Structures]: Performance Analysis and Design Aids—*simulation*; C.4 [Computer Systems Organization]: Performance of Systems; D.4.1 [Operating Systems]: Process Management; D.4.2 [Operating Systems]: Storage Management—*Allocation / deallocation strategies, main memory*; E.2 [Data]: Data Storage Representations

General Terms: Algorithms, Measurement, Performance

1. INTRODUCTION

Most general-purpose computer systems support both paged virtual memory [7] and caches [22]. Virtual-memory caches pages by translating virtual addresses to their corresponding real addresses (often by consulting a page

A preliminary version of the material presented here appeared as a technical report [16]. This research is supported in part by the National Science Foundation's Presidential Young Investigator and Computer and Computation Research Programs under grants MIPS-8957278 and CCR-8902536, A.T. & T. Bell Laboratories, Cray Research, and Digital Equipment Corporation. R. E. Kessler has been supported by a summer internship from Digital Equipment Corporation and by graduate fellowships from the National Science Foundation and the University of Wisconsin, Madison. He is employed currently at Cray Research, Inc.

Authors' addresses: R. E. Kessler, Cray Research, Inc., 900 Lowater Road, Chippewa Falls, WI 54729. email: kessler@cs.wisc.edu; M. D. Hill, University of Wisconsin, Computer Sciences Department, Madison, WI 53706.

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.

© 1992 ACM 0734-2071/92/1100-0338 \$01.50

ACM Transactions on Computer Systems, Vol. 10, No. 4, November 1992, Pages 338–359.

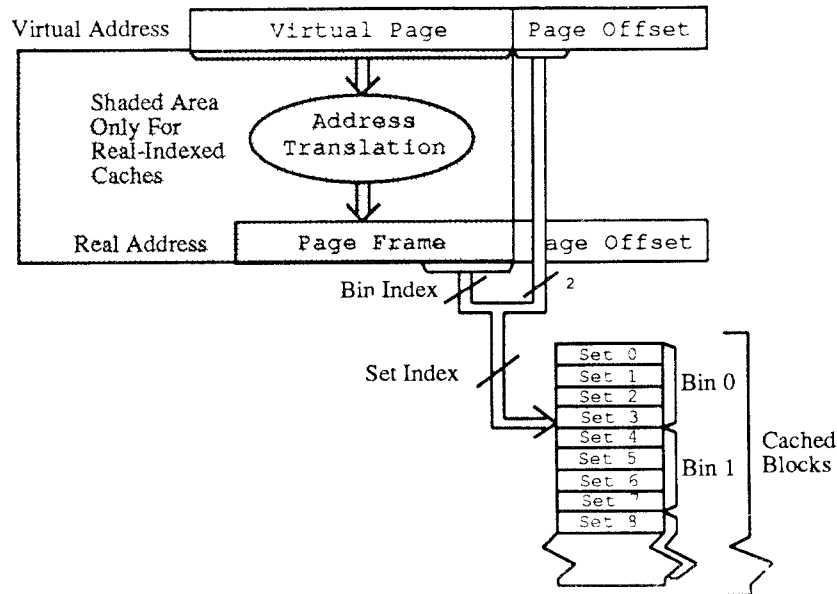


Fig. 1. Virtual memory and set-indexing interaction. This figure shows translation stages that index a virtually addressed memory location into a large set-associative CPU cache. The cache takes index bits directly from the virtual address for a virtually indexed cache, and it takes the bits from the real address for a real-indexed cache. Any of the A (A is the associativity) cache block frames in the set chosen by the *index* bits can contain the addressed memory location. This figure depicts bins that are four-cache sets. This paper only considers large real-indexed caches with bit selection indexing [22].

table). This translation from virtual pages to real page frames is usually fully associative: *any* page frame in the physical main memory can hold any page. In contrast, hardware caches usually have a more restrictive set-associative mapping that only allows a cache block to reside in one of the few frames in a set. A real-indexed (virtual-indexed) cache selects the set for an incoming reference by extracting index bits from the reference's real (virtual) address.

Figure 1 depicts virtual-address translation together with cache set indexing. A large real-indexed cache requires so many index bits that some of them must come from the address translation output.¹ We call the index bits coming from the address translation the *bin-index* and the group of sets selected by them a *bin*.² Address translation affects large real-indexed cache performance because page translation (implicitly) selects the cache bin that the cache blocks in the page reside in. In this study, we examine this interaction between page placement and cache performance. We focus on

¹ A large cache is defined to be one whose size divided by its associativity exceeds the page size. The associativity of the cache is the number of cache block frames in each set.

² Note that the *superset* of Goodman is different from a bin. Pieces of a superset span *all* bins.

large real-indexed caches because page placement has no effect on a cache that is either small or virtual indexed.

Two trends warrant an examination of the interaction between virtual-memory page mapping and large real-indexed caches. First, real-indexed caches are commonly implemented because they handle virtual-address synonyms (aliases) and virtual-address space changes (e.g., context switches) more simply than do virtual-indexed caches [22]. Real indexing is particularly appropriate for secondary caches (e.g., level-two caches: those that service the references from primary, or level-one, caches) in hierarchical CPU cache configurations, because the address translation cost of real indexing has a smaller performance impact [24, 26]. Second, typical cache sizes are growing, so large caches, even multimegabyte ones, are now feasible design alternatives [5]. Page mapping has a greater effect on these large caches because they have many more bins.

1.1 Previous Work

Few previous studies examine the page-mapping interaction with real-indexed cache performance because large real-indexed caches are a recent trend. Sites and Agarwal [21] compare the performance of virtual-indexed caches to real-indexed caches. They find that real-indexed caches perform worse than corresponding virtual-indexed caches unless there is frequent context switching.

Only a few systems optimize their page placements for cache performance. In the Sun 4/110, the operating system maps instruction pages and data pages [14] to even and odd page frames to partition the instructions and data in the 4/110's SCRAM cache [10]. MIPS uses a variant of *Page Coloring*, described in Section 3, to improve and stabilize its real-indexed cache performance [24]. With track-driven simulation, Bray et al [6] verify that Page Coloring reduces the mean and variance of the cache miss ratio, particularly for large caches. They also, model the effective main memory size reduction of a strict version of Page Coloring that assumes a page *must always* be placed in a frame of the correct color. This model does not apply to our work (or to MIPS' implementation) because it implies that page replacement may be required even when frames (of the wrong color) are available. Rather, we require page replacement *only* when there are no available page frames.

Finally, many others [8, 12, 17, 18, 23] have examined changing virtual-address reference patterns to improve cache performance. These techniques can be used in conjunction with our work, since our techniques alter real addresses, not virtual addresses.

1.2 Contributions of this Paper

Operating systems place (map) new pages into main memory by selecting from a pool of available page frames. This pool consists of page frames that have not been used recently and can be reallocated to new pages. The selection process is usually arbitrary; often the first available frame in the pool is reclaimed by the new page. We call this *random* (or naive) page

mapping because the operating system does not care *where* it places the new page. We introduce better page placement algorithms that we call *careful page-mapping* algorithms. They are small modifications to the existing operating system virtual-memory management software that improve *dynamic* cache performance with better *static* page bin placement decisions. When the pool of available page frames allows some mapping flexibility, rather than just choosing an arbitrary page (like the first available), a careful-mapping algorithm chooses one of the available page frames that best suits some simple static cache bin placement heuristics.

These careful-mapping algorithms are low overhead, particularly since they may only execute once when each page is mapped, while they may improve cache performance on each reference to the page. Furthermore, our careful-mapping algorithms need not impact the page fault rate by changing page replacement decisions, because they never force a page replacement when there are any page frames in the available page frame pool, and because the contents of the available page frame pool are determined solely by the operating system's page replacement algorithm, not by the careful-mapping algorithms. Our algorithms provide the opportunity for cache performance improvements when the operating system has more frames available for placement. When main memory is underutilized, careful-mapping algorithms are a cheap way to improve cache performance because it is easy to have a large available pool. When main memory is highly utilized, only a small pool will be available, so careful mapping will be less effective (and less needed because page swapping may dominate performance).

Section 2 motivates the static improvements that these careful-mapping algorithms rely on. It shows first how the operating system can improve static page bin placement. A simple static analysis then shows the potential careful-mapping static improvements: as many as 30% of the pages from an address space unnecessarily conflict in a direct-mapped cache when using random mapping. Also, this analysis correctly predicts that the largest gain from careful page mapping comes when the cache is direct mapped.

Section 3 describes Page Coloring and introduces several careful page-mapping algorithms. *Page Coloring* matches the real-indexed bin to its corresponding virtual-indexed bin. *Bin Hopping* places successively mapped pages in successive bins. *Best Bin* selects a page frame from the bin with the fewest previously allocated and most available page frames. *Hierarchical* is a tree-based variant of Best Bin that executes in logarithmic time and produces cache size-independent placement improvements (i.e., it improves static placement in many different cache sizes simultaneously).

Section 4 discusses some of the methods used to get the trace-driven simulation results of this paper. Section 5 uses the simulation results to measure the cache miss reduction from the careful-mapping algorithms. It shows that careful page mapping reduces the direct-mapped cache miss rate by 10–20% with our traces. This is about half of the improvement available from doubling either the size or the associativity of a direct-mapped cache. Section 5 also compares the performance of the different, careful page-mapping implementations for one multiprogrammed trace, showing that the

Page Coloring implementations perform poorly and that the Best Bin results in the fewest cache misses.

The implication of this research is that careful mapping can improve the performance of systems with large real-indexed caches. We believe that careful mapping is appropriate for future operating system implementations since it improves cache performance with small software changes and no hardware changes.

2. MOTIVATION FOR CAREFUL PAGE MAPPING

This section shows the potential static real-indexed page placement improvement of careful mapping over random mapping. Figure 2 depicts two placements of the code, data, and stock pages from an address space in the real-indexed cache; the left one might be a random placement of the pages, and the right one might be a careful placement of the same pages. The random placement puts many pages in the same cache bins; this competition is undesirable since it can cause more cache *conflict misses* [11]. Some of the randomly mapped pages could be placed in the unused bins to reduce conflict misses.

We measure the potential static careful mapping improvement quantitatively by counting the *page conflicts*, C , which is the number of pages in a bin beyond the associativity of the cache. If u pages land in a cache bin with associativity A , then we say $\max(0, u - A)$ conflicts occur.³ With a direct-mapped cache there are six conflicts in the random mapping of Figure 2, and there are four conflicts in the careful mapping. Similarly, with a two-way set-associative cache, there are three and zero conflicts.

We can calculate the average page conflicts from random mapping if we assume each page frame is independently and equally likely to hold a page.⁴ Let N be the cache size in pages and $B = N/A$ be the number of bins. A binomial distribution gives the probability that exactly u of U pages fall in one of the B bins:⁵

$$P(u \text{ in bin}) = \binom{U}{u} \left(\frac{1}{B}\right)^u \left(1 - \frac{1}{B}\right)^{U-u}. \quad (1)$$

Using Eq. 1, Eq. 2 calculates the expected conflicts resulting from a random placement of U pages across all cache bins:

$$C_{avg} = B \sum_{u=A+1}^U (u - A) \cdot P(u \text{ in bin}). \quad (2)$$

³ We use $u - A$ rather than u because we feel that it predicts more correctly the magnitude of the conflict. If we used u , then conflicts would be minimized by first filling the cache (placing A pages in each bin) and then placing all the rest of the pages in a single bin; it is unlikely that this mapping would minimize dynamic cache misses.

⁴ We assume an infinite main memory, but the results for a finite main memory are similar for the parameters we consider.

⁵ Both Thiebaut and Stone [25] and Agarwal et al. [1] use binomial distributions for similar analysis.

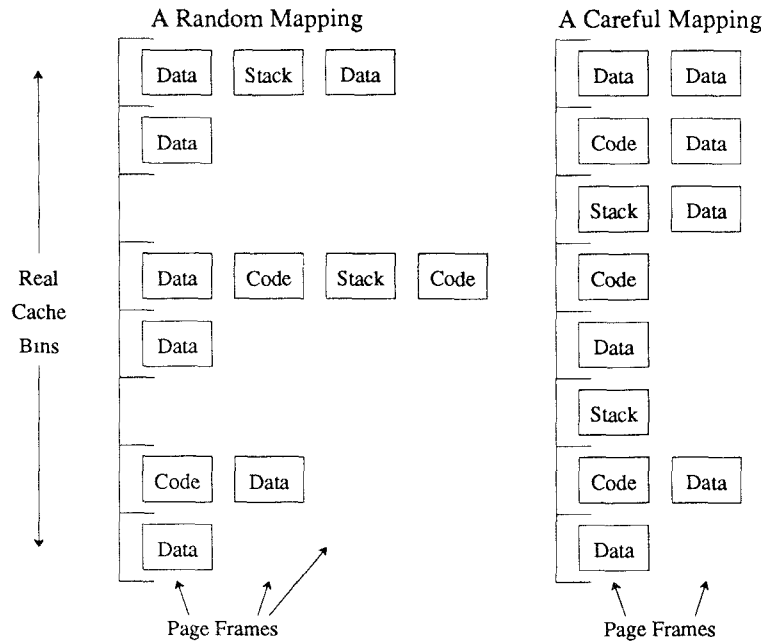


Fig. 2. Random and careful real-indexed bin placements. This figure shows two mappings of the code, data, and stack pages of an address space into a real-indexed cache. The left bin placement might result from random mapping, and the right one might be a careful mapping. The horizontally stacked pages are the (competing) ones that the virtual-to-real address translation indexes to the same bin.

The left graph in Figure 3 plots C_{avg} and its bounds for various address space sizes (U) with a direct-mapped cache. In the worst case, almost all the pages cause conflicts because they are placed in the same bin, so C_{max} is linear with the address space size. In the best case, the pages are evenly spread across the cache bins, so no conflicts occur until the cache is full of pages; and then each additional page adds a new conflict. C_{avg} remains closer to C_{min} than C_{max} , particularly for small and large values of U .

The right graph in Figure 3 illustrates the difference between C_{avg} and C_{min} more dramatically for direct-mapped, two-way, and four-way set-associative caches. It plots what we call the *mapping conflicts* ($C_{avg} - C_{min}$), which is the random conflicts less the minimum, and the potential static conflict savings of careful page mapping. It shows that the mapping conflicts are maximized at the point where the address space size equals the cache size. $C_{avg} - C_{min}$ is low for small address spaces since many collisions by few pages are unlikely. It is also low for large address spaces since the bins are likely to be highly utilized regardless of the page mapping, so there is little potential for improvement. This data shows that careful page mapping can reduce the static conflicts from random mapping; as many as 30% of the pages from an address space are unnecessarily in conflict with random mapping. It shows also that careful mapping is most useful for direct-mapped

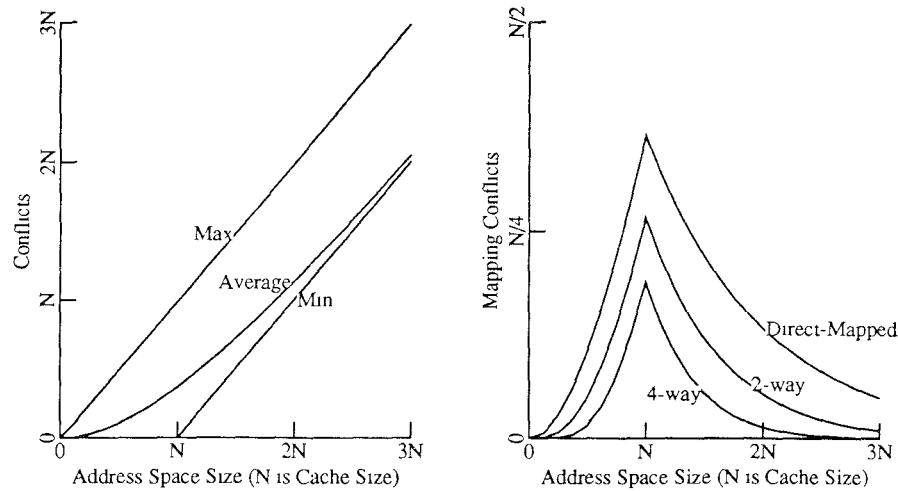


Fig. 3. Static cache placement conflicts. This figure shows the potential improvement of careful mapping by comparing the average number of static page conflicts in a random mapping to bounds. The left graph shows C_{avg} (as calculated by Eq. 2), $C_{max}(max(0, U - A))$, and $C_{min}(max(0, U - N))$ for various address space sizes (U) with a direct-mapped ($A = 1$) 64-bin cache ($N = B = 64$). The right graph shows the difference between the conflicts resulting from a naive mapping and the minimum conflicts, the mapping conflicts ($C_{avg} - C_{min}$), for different address space sizes and cache associatives ($A = 1, 2, 4, N = 64, B = N/A$).

caches because they have more static page conflicts. In the opposite extreme, there are never any mapping conflicts in a fully associative cache.

3. IMPLEMENTING CAREFUL MAPPING

This section presents four page-mapping implementations. A key aspect of all these algorithms is that they do not use any dynamic referencing information to improve page placements. Rather, they place pages into the bins of a large unified cache using solely static criteria.⁶ Static placement decisions are appropriate because dynamic information, such as cache references (misses) or page references (faults), may not be available to the page mapper. Section 5 shows that the static improvements of these algorithms eliminate many dynamic cache misses.

We consider systems in which multiple address spaces share the same main memory, so both intra-address-space and inter-address-space conflicts occur. Except for certain obvious circumstances, the careful-mapping algorithms primarily try to remove intra-address-space conflicts and only secondarily try to reduce inter-address-space conflicts. This is appropriate for

⁶ We do not consider mapping functions that statically improve the page placement in split caches. An example of split caches is split instruction and data caches [22]. With split caches, the mapping function may want to minimize the contention in both caches concurrently (or independently).

future high-performance processors because the memory within a *single* address space may be referenced many thousands or even millions of times before any inter-address-space conflicts cause any cache misses (only as a result of context switches).

In order to differentiate strictly between page replacement (making a frame available for reuse by placing it in the available pool) and page placement (choosing a frame from the pool), the careful-mapping algorithms always select frames that were previously in the available pool. Thus, page replacement needs to occur only when the pool is empty.⁷ Since the placement algorithms do not *require* an available pool with any specific characteristics and since many operating systems try to maintain a pool anyway (to smooth out peaks in memory demand, for example, [2]), we consider placement independently from how the available pool is replenished (either by page deallocation or replacement). The placement algorithms need not force *any* page replacement modifications, though clearly the size of the available pool can determine the page placement quality, and too large a pool can increase the page fault rate because of poorer page replacement decisions. Since the pool size may be limited, careful-mapping algorithms should preferably produce better placements with smaller pools.

3.1 Page Coloring

The simplest careful-mapping technique simply makes a real-indexed cache perform like a virtual-indexed cache by matching bits of the virtual page with bits from the real page frame number. MIPS has implemented a variant of this called Page Coloring [24]. In our simple implementation of Page Coloring, the page-mapping function simply examines the pool of available page frames to find one with the proper bit value(s). When it finds a match, it extracts the frame and places the page in it. When there is not a matching frame in the pool, our implementation chooses an arbitrary frame (the first available in the pool). Given the complexity of current memory systems the simplicity of Page Coloring is a virtue.

Page Coloring minimizes cache contention because sequential virtual pages do not conflict with each other in the cache. Mostly contiguous address spaces (the common case) will usually have little cache contention with Page Coloring. However, the exact match of virtual bits with real bits, as described above, can lead to excessive inter-address-space contention in a multiprogrammed workload because it places commonly used virtual addresses from different address spaces (like the stack, for instance) in the same bin. We consider two versions of Page Coloring: (1) an exact match of bin bits and (2) a match of bin bits exclusive-ored with the process' identifier (PID hashing). PID hashing makes it unlikely that the same virtual address from two different address spaces maps to the same bin, thereby reducing contention when some virtual addresses are frequently used (e.g., if the top of stack for

⁷ Bray et al. [6], assume that a page is always placed in a frame of the correct color, even if a page replacement must be initiated to get one.

all processes is at the same virtual address). It is our understanding that MIPS uses a variant of PID hashing [24].

3.2 Bin Hopping

Bin Hopping is another simple way to equalize cache bin utilization. It allocates sequentially mapped pages in sequential cache bins, irrespective of their virtual addresses. Bin Hopping reduces cache contention because it sequentially distributes the pages from an address space across different bins (until it must wrap around). It exploits temporal locality because the pages it maps close in time tend to be placed in different bins. If pages mapped together are also referenced together, this will be a useful property for reducing cache contention.

The Bin Hopping implementation simply traverses the bins until it finds one of them that has an available page frame. Then it remembers where it placed the last page, so the next traversal can start from there. It maintains a separate traversal path per address space to minimize contention within each address space separately. The simplicity of Bin Hopping, like Page Coloring, is its largest asset. An advantage of Bin Hopping over Page Coloring is that it is more flexible when there are few pages available. Instead of choosing an arbitrary page when one is not available in the desired bin, it skips to another bin that is almost as good as the previous one. Bin Hopping is able to do this because it remembers the last page placement for each address space.⁸

3.3 Best Bin

The strength of Page Coloring and Bin Hopping could also be considered their weakness. They both have simple implementations that require little (or no) storage, but their decision making is not very sophisticated. When there are few available page frames, Page Coloring and Bin Hopping may not produce a good mapping.

Best Bin uses the counter pair $\langle used, free \rangle$ associated with each bin to select a bin. *Used* is the number of previously placed pages from a given address space, and *free* is the number of page frames available (system wide) in this bin.⁹ Best Bin reduces contention by placing a page in the bin with the smallest *used* and largest *free* value. The choice of the smallest *used* value minimizes intra-address-space contention because all bins become equally used by each address space, and the secondary choice of the largest *free* value tends to minimize inter-address-space contention because large free values may imply bin underutilization.

⁸ We initialize this state to a random location before any mapping occurs.

⁹ While it is possible to determine these counts from the complete mapping of virtual pages to page frames, Best Bin can store them in a more readily available form for greater execution efficiency. The added storage required by Best Bin is the *used* counter per bin in each address space, since the space to store the single system-wide *free* counters is minimal. Even the *used* counter storage is modest when compared to the complete page-mapping information.

As an example of a choice made by Best Bin, consider a cache with four bins and $\langle used, free \rangle$ pairs of $\langle 0, 0 \rangle$, $\langle 1, 3 \rangle$, $\langle 1, 1 \rangle$, and $\langle 2, 4 \rangle$. Best Bin would select the $\langle 1, 3 \rangle$ bin; $\langle 0, 0 \rangle$ could not be selected because there are no available frames, even though it has the minimal *used* count; $\langle 1, 3 \rangle$ is better than $\langle 1, 1 \rangle$ because there are more available frames, and $\langle 1, 3 \rangle$ better than $\langle 2, 4 \rangle$ because its *used* count is smaller. After choosing it, Best Bin modifies the $\langle 1, 3 \rangle$ bin to $\langle 2, 2 \rangle$ to reflect the altered system state.

3.4 Hierarchical

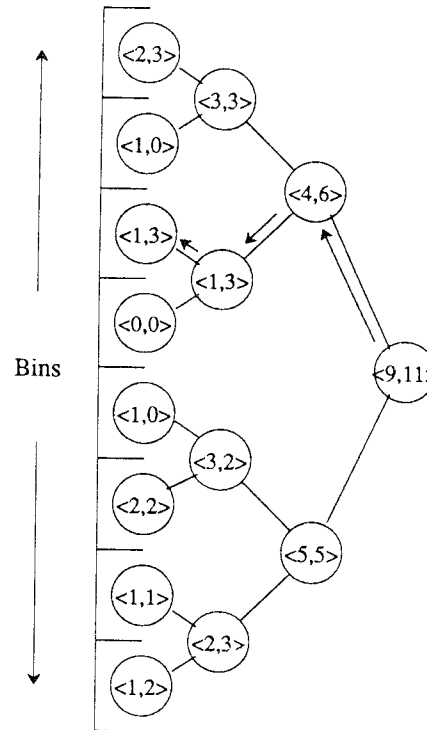
The execution time of Best Bin may be linear in the number of bins, which is too slow when there are too many bins. Alternatively, Hierarchical can choose a bin in logarithmic time with a tree whose leaves contain $\langle used, free \rangle$ values for each bin and where each interior node value contains the sum of the values of its children.¹⁰ Hierarchical selects a bin by beginning at the tree's root and traversing down to its leaf. At each interior node, Hierarchical chooses a traversal direction by comparing $\langle used, free \rangle$ values of the children, using the same priorities as Best Bin.

Figure 4 shows Hierarchical on a simple example. Traversal starts from the root of the bin tree. From there, Hierarchical decides, based on the $\langle 4, 6 \rangle$ and $\langle 5, 5 \rangle$ pairs corresponding to its subtrees. It chooses $\langle 4, 6 \rangle$ since it has less *used* pages. Next, Hierarchical chooses $\langle 1, 3 \rangle$ over $\langle 3, 3 \rangle$, again because there are less *used* pages. Finally, it chooses $\langle 1, 3 \rangle$ over $\langle 0, 0 \rangle$ because $\langle 1, 3 \rangle$ has available page frames. Bin tree traversal is then complete since a leaf node is reached. Note that Best Bin would choose the same bin in this example. To update the state to account for the page mapping, Hierarchical modifies the values of the traversed nodes to $\langle 2, 2 \rangle$, $\langle 2, 2 \rangle$, $\langle 5, 5 \rangle$, and $\langle 10, 10 \rangle$ from bottom to top, respectively.

An important feature of Hierarchical is its cache size independence. Hierarchical can simultaneously optimize page placement for many cache sizes if the tree is constructed so the children of each internal node have bin numbers with the most least-significant bits in common [15]. For example, the tree for eight bins has eight leaf nodes 000, 001, ..., 111. One level higher up in the tree, there are only four nodes: 00, 01, 10, 11. For Hierarchical to be size independent 00 should parent 000 and 100; 01 should parent 001 and 101; 10 should parent 010 and 110, and 11 should parent 011 and 111. That way, each tree level selects one more (high-order) bit, or equivalently, optimizes for the next larger cache. For example, size independence is extremely important with a hierarchy of large real-indexed caches, since the performance of *all* caches (e.g., both primary and secondary) can be improved with a *single* size-independent page placement. Size-independent page placement also has

¹⁰ Hierarchical actually uses a single system-wide free bin tree, but it needs a used bin tree with each address space, just like Best Bin needs a bin array per address space. The storage for the used bin tree need not be more than twice that required for a bin array, still a small amount compared to the complete page-mapping information. This study uses binary trees; higher branching factors also could be used, but they may not maintain size independence with all caches.

Fig. 4. Bin tree traversal. This figure shows an example of a bin choice when mapping a virtual page by $\langle used, free \rangle$ bin tree traversal with Hierarchical. The values of the tree nodes are given. The used values of the leaves show the number of pages already placed in the corresponding bin. The free values of the leaves show the number of page frames available for replacement in the corresponding bin.



a simulation advantage: many different-size caches can be simultaneously simulated.

4. TRACE-DRIVE SIMULATION METHODOLOGY

This section discusses some of the methods used for the simulation results shown in Section 5.

4.1 The Traces

The traces used in this study were collected at DEC Western Research Laboratory (WRL) [4, 5] on a DEC WRL Titan [19], a load/store (“RISC”) architecture. Each trace consists of the execution of three to six *billion* instructions of large workloads on a load/store architecture, including multi-programming but not operating system references. Each trace references from eight megabytes to over 100 megabytes of unique memory locations. The traces are sufficiently long to overcome the cold-start intervals of even the very large caches considered in this study; we present cold-start results, but the error cannot be more than a few percent for even the largest cache. We chose programs with large memory requirements because large application sizes will likely be more common as main memories of hundreds of megabytes become common.

The traces of the multiprogrammed workloads represent the actual execution interleaving of the processes on the traced system with switch intervals

that may be appropriate for future high-performance processors. The **Mult2** trace includes a series of compiles, a printed circuit board router, a VLSI design rule checker, and a series of simple UNIX commands, all executing in parallel (about 40 megabytes active at any time) with an average of 134,000 instructions executed between each process switch. The **Mult2.2** trace is the Mult2 workload with a switch interval of 214,000 instructions. The **Mult1** trace includes the processes in the Mult2 trace plus an execution of the system loader (the last phase of compilation) and a Scheme (Lisp variant) program (75 megabytes active) and has a switch interval of 138,000 instructions. The **Mult1.2** trace is the Mult1 workload with a switch interval of 195,000 instructions. **Sor** is a uniprocessor-successive overrelaxation algorithm that uses large sparse matrices (62 megabytes). **Tv** is a VLSI timing verifier that first builds a large data structure and then traverses it with little locality (96 megabytes). **Tree** is a Scheme program that searches a large tree data structure (64 megabytes). **Lin** is a power supply analyzer that uses sparse matrices (57 megabytes).

4.2 The Simulator

Previous trace-driven cache simulation results either assume small or virtual-indexed caches or that operating system policies are fixed and independent of the studied caches. Our simulator incorporates memory management policies of the operating system along with hardware cache management implementations because we examine the interaction of hardware and software memory system components.

We arbitrarily chose a global least recently used (LRU) page replacement policy for this study. The simulator maintains an exact ordering of the page frames from most recently used to least recently used (an LRU list, randomly initialized when simulation begins, to model a system that has been executing for some time). We assume a page placement policy can select any page frame from the pool of available page frames at the end of the LRU list. (The simulator does not place frames in the available pool when a process exits and when its pages are deallocated, only when they are least recently used). The size of the available page frame pool is constant throughout the simulations of this study, even at the beginning of a simulation when the entire main memory is unmapped. The default pool size is 4 megabytes of 16-kilobyte pages. The main memory size of 128 megabytes is large enough that the simulations required no process swapping, so we use the static execution interleaving captured on each trace to model context switching.

The simulated CPU cache configuration is a two-level hierarchical system that does not maintain multilevel inclusion [3]. The primary (level-one) caches are split direct-mapped write-allocate, write-back instruction and data caches of 32 kilobytes each with 32-byte blocks.¹¹ They are backed up by a 1-megabyte to 16-megabyte unified secondary (level-two) write-back cache

¹¹ On a write miss, space is allocated, and the cache block is read in (write-allocate). On a miss to the data cache that requires a write of a dirty block, the write-back is executed followed by the block read.

with a random replacement policy and 128-byte blocks. We simulate many different secondary cache configurations, and we focus on the large secondary caches since they are most likely to be real indexed and have the most to gain from careful page mapping.

4.3 Misses per Instruction

Our simulation tools report cache performance in *misses per instruction* (MPI). In this study, we only compare alternatives for the large secondary caches, so the MPI results are proportional to global miss ratio (cache misses divided by processor references [20]), and relative changes in MPI are equal to relative changes in global miss ratio. We do not use local miss ratio (cache misses divided by cache accesses [20]), because careful mapping can cause the local miss ratio of a secondary cache to increase even as system performance improves (by reducing primary cache misses and hence secondary cache accesses). Kaplan and Winder [13] provide good arguments for using misses per instruction.

5. TRACE-DRIVEN SIMULATION RESULTS

This section shows first that the static page placement improvements of careful page mapping (*Hierarchical*) eliminate many dynamic cache misses. Then it compares the alternative careful-mapping implementations using one of the multiprogrammed traces.

5.1 The Usefulness of Careful Page Mapping

We compare Hierarchical mapping to random mapping in this section because Hierarchical's size independence reduces simulation time and because Hierarchical performs well compared to the other careful-mapping schemes (as will be shown in Section 5.2). For 1-megabyte (top), 4-megabyte, and 16-megabyte caches (bottom), Figure 5 plots 100 million instruction averages of the Mult2.2 misses per instruction (MPI) for four different simulations, each with a different virtual-to-real page mapping. The cache performance is different for each simulation because they each place the same pages in different frames. Note that random mapping *does not* imply random page replacement; it means that the simulation places the page in the page frame at the bottom of the LRU list. Thus, in the random runs the simulator places the page in a random (arbitrary) bin that is determined by the page frame number at the bottom of the LRU list.

For the Mult2.2 trace, Figure 5 shows that the careful (Hierarchical) mappings have consistently lower MPI than random (since the solid lines are usually below the dotted lines). This clearly indicates that the static page placement improvements of careful mapping do reduce dynamic cache misses for this trace. Over most of the trace, the careful mapping MPI's are only a few percent better than random. But in addition to this small and constant improvement, there are bursts of cache misses from poor random page mapping that occasionally double or triple the MPI for the 100 million

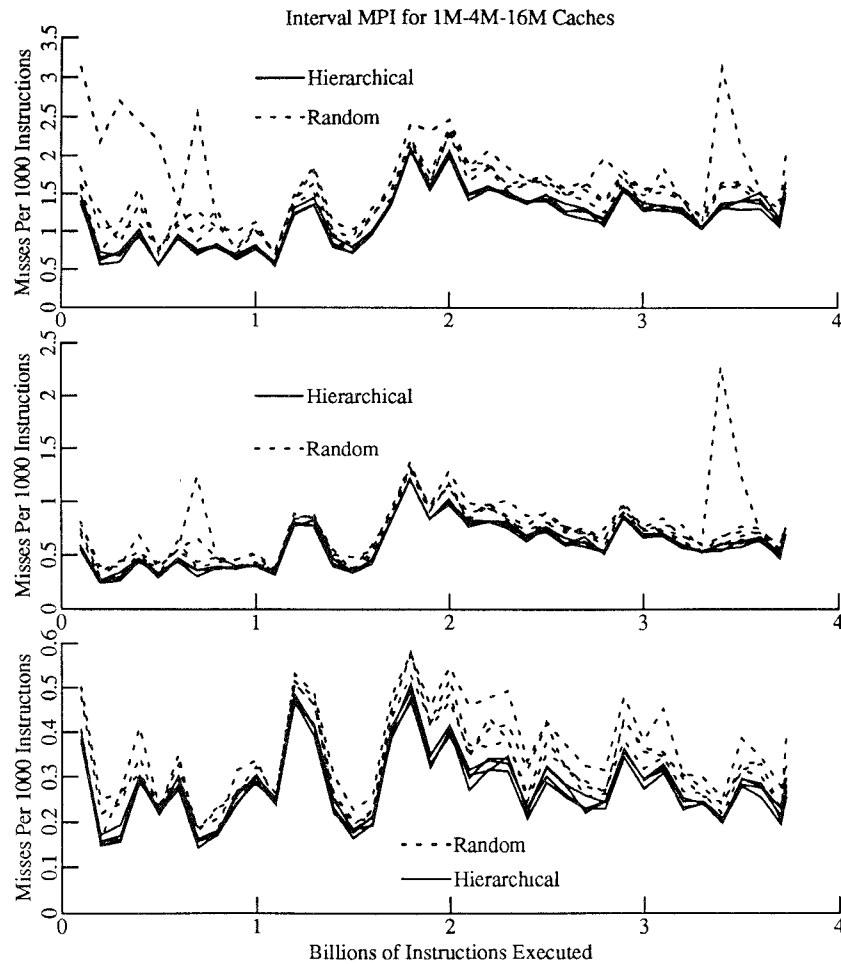


Fig. 5. Mult2.2 results for different page mappings. This figure shows the real-indexed MPI resulting from four different virtual-to-real mappings generated by Hierarchical (solid lines) and generated randomly (dotted lines) on the misses of three different secondary caches. It shows the average MPI for the previous 100 million instructions of the Mult2.2 trace for 1-megabyte (top), 4-megabyte (middle), and 16-megabyte (bottom) secondary direct-mapped caches. Note that the scales are different for each graph.

instruction intervals.¹² These bursts increase the mean improvement of careful mapping beyond only a few percent.

To get a better idea of the quantitative dynamic cache performance improvements of careful mapping, we wish to determine the average cache MPI

¹² Note that some bursts appear in the random mappings of both the 1-megabyte and 4-megabyte cache simulations. This occurs since the simulations of different cache sizes are of the same virtual-to-real page mapping, and the results are correlated. The results for each individual cache are not correlated, however.

when using random and Hierarchical mapping. To do this, we averaged the results from multiple simulation runs that use different page mappings. While the MPI of one simulation is an unbiased estimate of the true mean MPI, Figure 5 suggests that the MPI results from simulations can substantially vary, so we do not know how close one simulation result is to the true mean MPI. With a sample of size four, the variations across the simulations were sufficiently small that the four-sample mean could serve as a good estimate of the true mean (most simulations are within a few percent of the sample mean). It seems that each trace was long enough that the individual simulation runs already averaged out much of the variance in cache performance. Luckily, a sample of size four was sufficient, since each long-trace simulation takes days. With shorter traces, we might have needed many more simulations to obtain the same accuracy.¹³ Below, we use four-sample means to quantitatively evaluate the improvement of careful mapping for 1-megabyte, 4-megabyte, and 16-megabyte secondary caches with direct-mapped, two-way, and four-way set-associativity.

Figure 6 shows the direct-mapped Hierarchical miss reductions. This gives the relative amount that the four-sample Hierarchical MPI is lower than the four-sample random-mapping MPI. The figure shows that careful page mapping is useful for a variety of workloads. Careful mappings eliminate up to 55% of the direct-mapped cache misses from random mapping. On the average over all the traces, it reduces the random-mapping direct-mapped cache misses per instruction by about 10–20%. A similar look at the results for two- and four-way associativities suggests that careful mapping eliminates about 4–10% of the misses in a two-way set-associative cache with random replacement, and about 2–5% of the misses in a four-way set-associative cache with random replacement (on average). This verifies that the largest gain from careful mapping comes with a direct-mapped cache, as predicted by the simple static analysis of Section 2. Higher set-associativities can more easily reduce contention without static placement improvements because they spread contending blocks across the cache block frames in a set. Static improvements that avoid contention are more valuable with direct mapping because direct-mapped caches cannot handle contention for the same set (block frame).

The direct-mapped Hierarchical miss reduction is positive in all cases except with *Lin*, where the mean MPI slightly increases with careful mapping. The multiprogrammed trace behaviors are all similar, except for *Mult1.2*, where one bad page placement in one simulation biased the results and cut the direct-mapped Hierarchical reduction in about half (for the 1-megabyte and 4-megabyte caches). Careful mapping does not help *Tv* much because most of its misses occur while it traverses a large data structure with little locality; it cannot eliminate many misses during this traversal because the memory locations tend to be randomly accessed, and the *Tv* address space

¹³ The statistical significance of these and many more results can be checked by consulting Kessler's thesis [15]. In most cases, our results are accurate to within $\pm 5\%$.

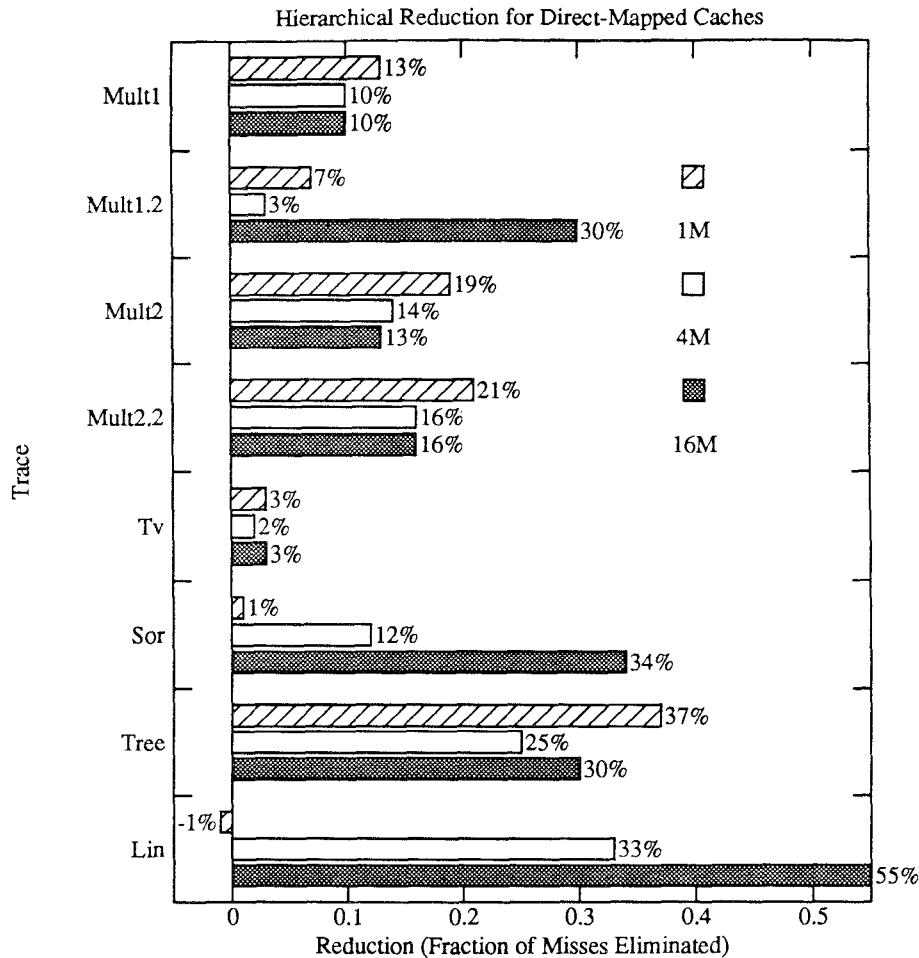


Fig. 6. Direct-mapped hierarchical miss reductions for all traces. This figure shows, for different direct-mapped secondary caches and traces, the Hierarchical reduction, which is the fraction of the misses (produced from random mapping) that Hierarchical could eliminate.

is too large for much static placement improvement. (Figure 3 shows how there is little room for static mapping improvements when the active address space size is large relative to the cache size.) The improvements of both Sor and Lin relate closely to the benefits of associativity. There is little careful mapping or associativity improvement with the 1-megabyte cache, but both careful mapping and associativity improve performance greatly for the 16-megabyte cache because some arrays fit fully in the cache. Though the results in Figure 6 vary greatly for the different traces, the rest of this study examines Mult2.2 in detail because of its “average” behavior; the Mult2.2 Hierarchical reductions are close to the average of all the traces.

Figure 7 shows the Mult2.2 results in another graphical way to put the careful-mapping miss reduction in its proper perspective. By interpolating

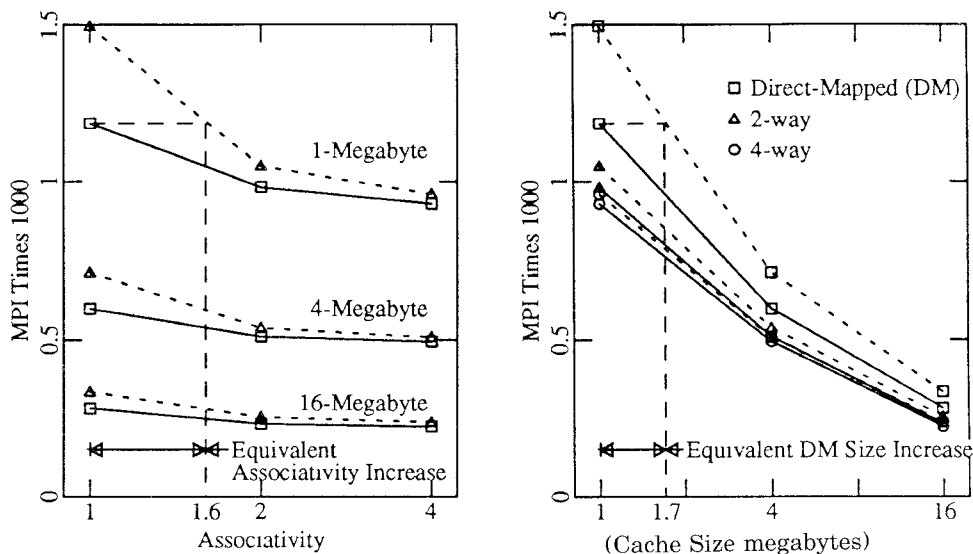


Fig. 7. Mult2.2 careful mapping improvement This figure shows the (mean) Mult2.2 simulation data. The misses per thousand instruction is plotted versus the cache associativity and size for random (dotted) page mappings and Hierarchical (solid) page mappings. The left graph shows MPI versus associativity for the different cache sizes. The right graph shows MPI versus cache size for different associativities (square = direct-mapped, triangle = two-way, and circle = four-way).

between the simulation results for different cache sizes and associativities, it estimates the effect of software careful-mapping relative to hardware cache design changes. The dashed line in the left graph finds the equivalent associativity change of Hierarchical with a 1-megabyte direct-mapped cache. Careful mapping reduces misses by more than half as much as a random-mapping associativity increase does; it is equivalent to an associativity increase from direct-mapped to 1.6-way with Mult2.2. For a direct-mapped cache, the dashed line in the right graph compares the careful-mapping miss reduction (with a 1-megabyte cache) to the reduction from a randomly mapped cache size increase. Careful mapping also removes nearly as many misses as a direct-mapped cache size doubling does; it makes a 1-megabyte direct-mapped cache act like a 1.7-megabyte direct-mapped cache for Mult2.2.

Most importantly, as we have shown, careful page mapping reduces mean MPI (across different page mappings) by avoiding poor page mappings. But a secondarily useful trait of careful mapping is a reduction in the MPI variance for different runs of a trace. Variance in cache behavior is undesirable because it increases the execution time variance of a workload. Figure 5 shows that the Mult2.2 mappings produced carefully have less MPI variations than the random mappings because all the solid (Hierarchical) lines are much closer together than the dotted (random) lines. The different careful-mapping simulations tend to follow an outline that the workload dictates, but

Table I. Careful-Mapping Comparison

Secondary Cache $MPI \times 1000$ (Mult2.2)		
Heuristic	Large Pool	Small Pool
Random	0.71 (19%)	0.71 (8%)
Page Coloring (equal)	0.69 (15%)	0.85 (30%)
Page Coloring (hash)	0.64 (8%)	0.71 (7%)
Bin Hopping	0.60 (0%)	0.67 (1%)
Best Bin	0.60 (0%)	0.61 (-8%)
Hierarchical	0.60 (0%)	0.66 (0%)

This table compares results of the different mapping heuristics outlined in Section 3 for different page frame pool sizes for the Mult2.2 trace. It shows the four-sample mean MPI for four megabytes of available page frames (Large Pool) and 256 kilobytes of available page frames (Small Pool). In parenthesis, it also gives the relative difference from the Hierarchical results for each pool size. The results shown are for a four-megabyte direct-mapped cache.

the random-mapping simulations show the bursts that are caused by random contention for cache locations. Poor random mappings into the cache, not the workload references themselves, cause many of the random-mapping misses. We found that when Hierarchical removes the most misses, it also tends to eliminate the most variance. This is not surprising since most of the MPI reduction comes from eliminating bursts of misses that increase both the mean and variance of the MPI. Hierarchical had a lower four-sample MPI variance than random mapping for 57 out of the 72 cases we examined. With 14 out of the 24 direct-mapped caches, the random-mapping four-sample variance is more than five times larger than Hierarchical. This strongly suggests that Hierarchical will yield more predictable execution times than random mapping.

5.2 Careful-Mapping Alternatives

In this section, we compare the Heuristics from Section 3 using the Mult2.2 trace. The size of the available page frame pool has an important effect on the comparison. The implementations that produce the fewest misses with the least available page frames are preferred because they are adaptable to the most situations. For the Hierarchical results given so far in this study the pool is 4 megabytes. As an experiment to learn the impact of the available pool size, Table I compares the 4-megabyte direct-mapped Mult2.2 MPI of the different heuristics (four-sample mean) with a 4-megabyte pool of available pages (large) and a 256-kilobyte pool (small).¹⁴ On the average, there was one available page frame in each bin with the large pool, and 0.0625 available page frames in each bin with the 256-kilobyte pool.

Equality Page Coloring performs poorly in the Mult2.2 comparison presented in Table I. Its particularly high mean MPI for the small pool comes

¹⁴ Most of these results are accurate to within ± 0.02 or less.

from the exceptionally high MPI of two (out of the four) simulations. This is a symptom of the basic problem with equality Page Coloring: the placement of commonly used virtual pages from different address spaces in the same cache bins often leads to many unnecessary cache misses. Probably because of this problem, MIPS does not use the strict equality match Page Coloring simulated here. The PID-hashed Page Coloring implementation performs much better for Mult2.2, yet it still performs poorly compared to the other careful-mapping schemes. The limited flexibility of the Page Coloring implementation is the cause of this poor performance: if Page Coloring does not find an available page frame in the matching bin, it abandons the search and chooses any arbitrary page frame. With the small pool, an exact match is unlikely (probability 0.0625), so the placement is almost random. The other careful-mapping schemes are more adaptable to situations with few available page frames.

The data in Table I shows that Best Bin produces superior placements because a smaller pool does not significantly affect its performance, while it causes the four simulations with the other heuristics to have more cache misses. With the 4-megabyte pool of available page frames in Table I, Bin Hopping, Best Bin, and Hierarchical perform nearly equivalently. Each of them finds it easy to reduce cache contention when there is one available page frame in each bin on average. With 0.0625 available page frames per bin Best Bin produces the fewest misses. This shows that Best Bin is the most flexible across all pool sizes and that it can improve Mult2.2 cache performance under a wide variety of page availability conditions.

We show results only for a 4-megabyte direct-mapped cache in Table I. Also, we found when using a single placement to simulate many cache sizes that Hierarchical eliminated the most misses over the range of sizes. This verifies the size-independent advantage of Hierarchical.

6. CONCLUSIONS

When a real-indexed cache is large, the virtual-to-real page-mapping function and the cache-mapping function interact to determine cache performance. This gives the page-mapping policy the opportunity to improve cache performance by carefully choosing the virtual-to-real page-mapping of an address space. This paper shows the usefulness of careful virtual-to-real page-mapping policies, describes several practical careful-mapping implementations, and shows their resulting cache miss reduction using trace-driven simulation. It is worthwhile to expend a little effort to map a page carefully since it can improve memory system performance on all following accesses to the page.

We feel that careful virtual-to-real page-mapping is a useful, low-cost technique to improve the performance of large, real-indexed CPU caches in many circumstances. We introduced careful-mapping algorithms that require no page replacement policy changes: they never force a page replacement when there are page frames in the available page frame pool, and the contents of the available page frame pool are determined solely by the operating system's page replacement algorithm. Our algorithms provide the opportunity for cache performance improvements when the operating sys-

tem makes more pages available for placement. When main memory is underutilized, careful-mapping algorithms are a cheap way to improve cache performance because it is easy to maintain a large available pool. When main memory is more highly utilized only a small pool will be available, so careful mapping will be less effective (and less needed because page swapping may dominate performance). An increase in the available page frame pool size gives the careful-mapping algorithms more freedom, perhaps with little or no corresponding page fault rate increase since a modest-sized pool provides sufficient flexibility. (Unfortunately, our traces were not long enough to quantify the relationship between page fault rate and pool size.)

This paper shows that static cache bin placement optimizations can produce dynamic cache performance improvements. A simple static analysis shows that random page placement causes many of the pages from an address space to be unnecessarily in conflict, so there is potential for static mapping improvements with careful mapping. This analysis correctly predicts that conflicts are less of a problem in caches of higher associativity.

This paper describes Page Coloring and introduces several other practical careful page-mapping heuristics that improve the cache bin placement of pages based solely on static information. These algorithms equalize cache bin utilization and eliminate much unnecessary cache contention. They range from the simple Page Coloring and Bin Hopping to the more sophisticated Best Bin and Hierarchical. This study quantitatively compares the MPI improvements of the heuristics with trace-driven simulations of a multiprogrammed trace (Mult2.2). This comparison shows that Bin Hopping, Best Bin, and Hierarchical perform well. Under this comparison, equality match Page Coloring performed poorly because it maps frequently used areas of the virtual-memory space from different address spaces to the same cache bin. A PID-hashed version of Page Coloring that fixes this problem still does not compare favorably to the other careful-mapping schemes. We found that Best Bin produces the best page placements with the smallest pool of available page frames; it is more adaptable across the range of possible page availability conditions. Hierarchical has good computational efficiency and cache size-independent mappings; these properties are desirable, particularly for hierarchies of real-indexed caches.

This paper presents simulation results from several workloads to measure the improved MPI mean and variance of careful page mapping (Hierarchical). The improvement depends on the workload, but with our traces we found that careful mapping reduced the direct-mapped cache misses per instruction by about 10–20%. Thus, at no hardware cost, software careful page mapping gets about half the improvement of either doubling the cache size or associativity.

ACKNOWLEDGMENTS

We would like to thank The Western Research Laboratory of Digital Equipment Corporation, especially Anita Borg and David Wall, for the traces used in this study. Joel Bartlett, Renato De Leone, Jeremy Dion, Norm Jouppi,

Bob Mayo, and Don Stark were a tremendous help in providing traceable applications. Paul Vixie and Colleen Hawk were a great help in storing the traces. Jim Goodman, Viranjit Madan, Guri Sohi, and David Wood from the Wisconsin Architecture Group provided valuable feedback. Paul Beebe and the Systems Lab were able to satisfy our *enormous* computing needs. Mike Litzkow and Miron Livny were a great help in adapting Condor to the requirements of these simulations. Special thanks to Harold Stone for his useful comments, and the referees for greatly improving the presentation of this paper.

REFERENCES

1. AGARWAL, A., HOROWITZ, M., AND HENNESSY, J. An analytical cache model. *ACM Trans. Comput. Syst.* 7, 2 (May 1989), 184–215.
2. BABAOLU, O., AND JOY, W. Converting a swap-based system to do paging in an architecture lacking page-referenced bits. In *Proceedings of the 8th Symposium on Operating System Principles*. 1981, pp. 78–86.
3. BAER, J., AND WANG, W. On the inclusion properties for multi-level cache hierarchies. In *Proceedings of the 15th Annual International Symposium on Computer Architecture*. 1988, pp. 73–80.
4. BORG, A., KESSLER, R. E., LAZANA, G., AND WALL, D. W. Long address traces from RISC machines: Generation and analysis. Res. Rep. 89/14, Western Research Laboratory, Digital Equipment Corporation, Palo Alto, Calif., 1989.
5. BORG, A., KESSLER, R. E., AND WALL, D. W. Generation and analysis of very long address traces. In *Proceedings of the 17th Annual International Symposium on Computer Architecture*. 1990, pp. 270–279.
6. BRAY, B. K., LYNCH, W. L., AND FLYNN, M. J. Page allocation to reduce access time of physical caches. Tech. Rep. CSL-TR-90-454, Computer Systems Laboratory, Stanford Univ., Stanford, Calif., 1990.
7. DENNING, P. J. The working set model for program behavior. *Commun. ACM* 11, 5 (May 1968), 323–333.
8. FERRARI, D. The improvement of program behavior. *IEEE Comput.* (Nov. 1976), 39–47.
9. GOODMAN, J. R. Coherency for multiprocessor virtual address caches. In *Proceedings of the 2nd International Conference on Architectural Support for Programming Languages and Operating Systems* (Oct.). 1987, pp. 72–81.
10. GOODMAN, J. R., AND CHIANG, M. The use of static column RAM as a memory hierarchy. In *Proceedings of the 11th Annual International Symposium on Computer Architecture*. 1984, pp. 167–174.
11. HILL, M. D., AND SMITH, A. J. Evaluating associativity in CPU caches. *IEEE Trans. Comput.* 38, 12 (Dec. 1989), 1612–1630.
12. HWU, W. W., AND CHANG, P. P. Achieving high instruction cache performance with an optimizing compiler. In *Proceedings of the 16th International Symposium on Computer Architecture*. 1989, pp. 242–251.
13. KAPLAN, K. R., AND WINDER, R. O. Cache-based computer systems. *IEEE Comput.* 6, 3 (Mar. 1973), 30–36.
14. KELLY, E. Personal communications. 1990.
15. KESSLER, R. E. Analysis of multi-megabyte secondary CPU cache memories. Ph.D. thesis, Computer Sciences Tech. Rep. #1032, Univ. of Wisconsin—Madison, 1991.
16. KESSLER, R. E., AND HILL, M. D. Miss reduction in large, real-indexed caches. Computer Sciences Tech. Rep. #940, Univ. of Wisconsin—Madison, 1990.
17. LAM, M. S., ROTHBERG, E. E., AND WOLF, M. E. The cache performance and optimizations of blocked algorithms. In *Proceedings of the 4th International Conference on Architectural Support for Programming Languages and Operating Systems* (Apr.). 1991, pp. 63–74.

18. MCFARLING, S. Program optimization for instruction caches. In *Proceedings of the 3rd International Conference on Architectural Support for Programming Languages and Operating Systems* (Apr.). 1989, pp. 183–191.
19. NIELSEN, M. J. K. Titan system manual. Res. Rep. 86/1, Western Research Laboratory, Digital Equipment Corporation, Palo Alto, Calif., 1986.
20. PRZYBYLSKI, S., HOROWITZ, M., AND HENNESSY, J. Characteristics of performance-optimal multi-level cache hierarchies. In *Proceedings of the 16th Annual International Symposium on Computer Architecture*. 1989, pp. 114–121.
21. SITES, R. L., AND AGARWAL, A. Multiprocessor cache analysis using ATUM. In *Proceedings of the 14th Annual International Symposium on Computer Architecture*. 1988, pp. 186–195.
22. SMITH, A. J. Cache memories. *Comput. Surv.* 14, 3 (Sept. 1982), 473–530.
23. STAMOS, J. W. Static grouping of small objects to enhance performance of a paged virtual memory. *ACM Trans. Comput. Syst.* 2, 2 (May 1984), 155–180.
24. TAYLOR, G., DAVIES, P., AND FARMWALD, M. The TLB slice—A low-cost high-speed address translation mechanism. In *Proceedings of the 17th Annual International Symposium on Computer Architecture*. 1990, pp. 355–363.
25. THIEBAUT, D., AND STONE, H. S. Footprints in the cache. *ACM Trans. Comp. Syst.* 5, 4 (Nov. 1987), 305–329.
26. WANG, W., BAER, J., AND LEVY, H. M. Organization and performance of a two-level virtual-real cache hierarchy. In *Proceedings of the 16th Annual International Symposium on Computer Architecture*. 1989, pp. 140–148.

Received May 1990; revised July 1991 and February 1992; accepted June 1992