

Specification of the PUMA memory management design

Bruce Jacob and Trevor Mudge

Advanced Computer Architecture Lab
EECS Department, University of Michigan
{bj,tm}@eecs.umich.edu

Tech Report CSE-TR-314-96
August, 1996

In this report we specify the memory management design of a 1GHz PowerPC implementation in which a simple design is a prerequisite for a high clock rate and short design cycle. The scheme involves no translation hardware such as a translation lookaside buffer or a page-table-walking state machine. However, it is just as efficient as hardware-managed address translation and is much more flexible. Modern operating systems such as Mach charge between 0.16 and 0.28 CPI for address translation on systems with TLBs. PUMA's software-managed address translation exacts an overhead of 0.03 CPI. Mechanisms to support such features as shared memory, superpages, sub-page protection, and sparse address spaces can be defined completely in software, allowing much more flexibility than in hardware-defined mechanisms.

Our software design combines the virtual caches with a PowerPC-like segmentation mechanism; it maps user 32-bit addresses onto a 44-bit segmented virtual space. We use a global page table to map this entire global virtual space. There are no individual per-process page tables; all process page tables map directly onto the global table and when processes share memory they also share a portion of the global page table. The following benefits are derived from the organization: (a) virtual cache consistency management can be eliminated, (b) the page table space requirements can be cut in half by eliminating the need to replicate page table entries for shared pages, and (c) the virtual memory system can be made less complex because it does not have to deal with the virtual-cache synonym problem.

1 Introduction

In this report we describe the hardware and software components of the PUMA virtual memory system, a memory management design that stays within an acceptable performance overhead and that does not require complex hardware. It places few constraints on the operating system but still provides all the features of systems with more hardware support. The hardware design is an implementation of *software-managed address translation*, or *softvm* for short [32]. It dispenses with hardware such as the translation lookaside buffers (TLBs) found in every modern microarchitecture and

the page-table-walking state machines found in x86 and PowerPC architectures. It uses a virtual cache hierarchy (8KB/8KB split L1 cache and 512KB/512KB split L2 cache—these specs current as of fall 1996) and takes an interrupt when a reference misses in the L2 cache. There is no table-walking hardware or TLB. The software-handled cache miss is similar to that of the VMP multiprocessor [11, 12, 13], except that VMP used the mechanism to explore cache coherence in a multiprocessor, while we use it to simplify memory management hardware in a uniprocessor. The design also resembles the in-cache address translation mechanism of SPUR [26, 44, 57] in its lack of TLBs, but takes the design one step further by eliminating table-walking hardware.

A *softvm* design performs several times better than a hardware-oriented design involving a TLB. It also supports common operating systems features such as address space protection, fine-grained protection, sparse address spaces, and superpages. Compared to more orthodox designs, it reduces hardware complexity without requiring unduly complex software. It has two primary components: a virtually indexed, virtually tagged cache hierarchy with a writeback cache at the lowest level (L2, for example), and a software-managed cache miss at the lowest level. Virtual caches do not require address translation when requested data is found in the cache, and so obviate the need for a TLB. A miss in the L2 cache invokes the operating system's memory manager, allowing the operating system to implement any type of page table, protection scheme, or replacement policy, as well as a software-defined page size. The migration of address-translation support from hardware to software increases flexibility significantly.

Virtual caches are an important part of the PUMA organization. They allow faster processing in the common case because they do not require address translation when requested data is found in the caches. However, they have not been used in many architectures despite their apparent simplicity because they have several potential pitfalls that need careful management [22, 30, 55]. We chose a virtual cache organization to meet the speed requirements of a high clock-rate pro-

cessor. We discovered that the segmented memory-management architecture of the PowerPC works extremely well with a virtual cache organization and an appropriate virtual memory organization, eliminating the need for virtual-cache management and allowing the operating system to minimize the space requirements for the page table. Management of the virtual cache can be avoided entirely if sharing is implemented through the global segmented space. This gives the same benefits as single address-space operating systems (SASOS): if virtual-address aliasing (allowing processes to use different virtual addresses for the same physical data) is eliminated, then so is the virtual-cache *synonym problem* [22]. Thus, consistency management of the virtual cache can be eliminated by a simple operating-system organization. The advantage of a segmented approach (as opposed to a SASOS approach) is that by mapping virtual addresses to physical addresses in two steps, a segmented architecture divides virtual aliasing and the synonym problem into two orthogonal issues. Whereas they are linked in traditional architectures, they are unrelated in a segmented architecture; thus applications can map physical memory at multiple locations within their address spaces—they can use virtual address aliasing—without creating a synonym problem in the virtual cache.

The PUMA software organization takes advantage of the PowerPC segmentation mechanism to eliminate virtual-cache synonyms and thus virtual-cache management. It shares memory through global addresses but the segmentation mechanism separates the structure of the global space from the process address space, so a process can attach a shared segment at any segment-aligned location in its address space, or even at multiple locations if desired. A global page table maps the global address space, and process page tables map directly onto the global table. There is therefore no per-process allocation of page tables; if two processes share memory they also share portions of the global page table. Thus, a process page table appears to be a contiguous, linear array of page table entries (PTEs), but it is in fact a disjunct set of pages in the global space, and therefore the organization is called a *disjunct page table*.

This *disjunct page table* is specific to a segmented architecture. With a typical degree of sharing in a system, the organization requires half the space of a normal page table. It eliminates the virtual-cache synonym problem by using global addresses for shared data; therefore the same virtual address is used for the same data and no synonym problem can occur. Without synonym problems, there is no need for management of the virtual cache, speeding up the system and simplifying its design immensely. The segmentation mechanism still supports virtual-address aliasing—the ability for a

process to use different virtual addresses (or even multiple virtual addresses) for the same physical data.

In this report, we present the details of our memory management system and describe it in the context of today's requirements for a virtual memory system.

2 Memory system requirements

There is a core set of functional mechanisms associated with memory management that computer users have come to expect, and that we wish to support in the PUMA design. They are found in nearly every modern microarchitecture and operating system (e.g., UNIX [3], Windows NT [15], OS/2 [16], 4.3 BSD [36], DEC Alpha [17, 47], MIPS [23, 33], PA-RISC [25], PowerPC [29, 40], Pentium [31], and SPARC [53]), and include the following:

Address space protection. User-level applications should not have unrestricted access to the data of other applications or the operating system. A common hardware assist uses address space identifiers (ASIDs), which extend virtual addresses and distinguish them from addresses generated by different processes. Alternatively, protection can be provided by software means [5, 19, 51].

Shared memory. Shared memory allows multiple processes to reference the same physical data through (potentially) different virtual addresses. Space requirements can be reduced by sharing code between processes. Using shared memory for communication avoids the data-copying of traditional message-passing schemes. Since a system call is typically an order of magnitude faster than copying a page of data, many researchers have investigated zero-copy schemes, in which the operating system unmaps pages from the sender's address space and re-maps them into the receiver's address space [18, 20, 37].

Large address spaces. Applications require increasingly large virtual spaces; industry has responded with 64-bit machines. However, a large address space does not imply a large address: large addresses are simply one way to implement large address spaces. Another is to provide each process a 4GB window into a larger global virtual address space, the approach used by the PA-RISC 1.X and 32-bit PowerPC architectures [25, 40].

Fine-grained protection. Fine-grained protection marks objects as read-only, read-write, execute-only, etc. The granularity is usually a page, though a larger or smaller granularity is sometimes desirable. Many systems have used protection to implement

various memory-system support functions, from copy-on-write to garbage collection to distributed shared virtual memory [2].

Sparse address spaces. Dynamically loaded shared libraries and multithreaded processes are becoming commonplace, but these designs require support for sparse address spaces. In contrast, 4.3BSD Unix [36] had an address space composed of two continuous regions. This design allowed the user page tables to occupy as little space as possible. Saving space was important, given that the original implementation did not allow page tables to be paged. A sparse address space has numerous holes in it, which would leave multiple holes within a page table; thus the wired-down, linearly-indexed page table of 4.3BSD would not be practical, as it would require 4MB of physical memory to map a 32-bit address space, regardless of how much or little of the virtual space is actually used.

Superpages. Some structures must be mapped for virtual access, yet are very large. The numerous page table entries (PTEs) required to map them flood the TLB and crowd out other entries. Systems have addressed this problem with “blocks” or “superpages”—multiples of the page size mapped by a single TLB entry. For example, the Pentium and MIPS R4000 allow mappings for superpages to reside in the TLB alongside normal mappings, and the PowerPC defines a Block TLB to be accessed in parallel with the normal TLB. Several studies have shown significant performance gains for reducing the number of TLB entries to cover the current working set [35, 48, 50].

Direct memory access. Direct memory access (DMA) allows asynchronous copying of data from I/O devices directly to main memory. It is difficult to implement with virtual caches, as the I/O space is usually physically mapped. The I/O controller has no access to the virtual-physical mappings, and so cannot tell when a transaction should first invalidate data in the processor cache. A simple solution performs DMA transfers only to uncached physical memory, but this could reduce performance by requiring the processor to go to main memory too often.

We will discuss our memory management design in the context of these features.

3 Background and previous work

Address translation is the mechanism by which the operating system provides virtual address spaces to

user-level applications. The operating system maintains a set of mappings from per-process virtual spaces to the system’s physical memory. Addresses are usually mapped at a *page* granularity—typically several kilobytes. The mappings are organized in a *page table*, and for performance reasons most hardware systems provide a *translation lookaside buffer* (TLB) that caches parts of the page table. When a process performs a load or store to a virtual address, the hardware translates this to a physical address using the mapping information in the TLB. If the mapping is not found in the TLB, it must be retrieved from the page table and loaded into the TLB before processing can continue.

3.1 Problems with virtual caches

Virtual caches complicate support for virtual-address aliasing and protection-bit modification. Aliasing can give rise to the *synonym problem* when memory is shared at different virtual addresses [22], and this has been shown to cause significant overhead [55]; protection-bit modification is used to implement such features as copy-on-write [1, 43], and can also cause significant overhead when used frequently.

The synonym problem has been solved in hardware using schemes such as dual tag sets [22] or back-pointers [52], but these require complex control logic that can impede high clock rates. Synonyms can be avoided by setting policy in the operating system—for example, OS/2 requires all shared segments to be located at identical virtual addresses in all processes so that processes use the same address for the same data [16]. SunOS requires shared pages to be aligned in virtual space on extremely large boundaries (at least the size of the largest cache) so that aliases will map to the same cache line [10, 24]¹. Single address space operating systems such as Opal [7, 8] or Psyche [46] solve the problem by eliminating the need for virtual-address aliasing entirely. In a single address space all shared data is referenced through global addresses; as in OS/2, this allows pointers to be shared freely across process boundaries.

Protection-bit modification in virtual caches can also be problematic. A virtual cache allows one to “lazily” access the TLB only on a cache miss; if so, protection bits must be stored with each cache line or in an associated page-protection structure accessed every cycle, or else protection is ignored. When one replicates protection bits for a page across several cache lines, changing the page’s protection can be costly. Obvious but expensive

1. Note that the SunOS scheme only solves the problem for direct-mapped virtual caches or set-associative virtual caches with physical tags; shared data can still exist in two different blocks of the same set in an associative, virtually-indexed, virtually-tagged cache.

solutions include flushing the entire cache or sweeping through the entire cache and modifying the affected lines.

3.2 Fragmentation of the virtual space

Most of the solutions to the synonym problem described above address the consistency problem by limiting the choices where a process can map a physical page in its virtual space. In some cases, the number of choices is reduced to one; the page is mapped at one globally unique location or it is not mapped at all. While this would seem to be a simple and elegant way to solve the virtual cache consistency problem, it creates another headache for operating systems, namely that of fragmentation.

When a global shared region is garbage-collected, the region cannot help but become fragmented. The problem is that whereas de-fragmentation of physical memory or disk space is as simple as copying pages or blocks, virtual pages cannot be relocated by simple copying because they are location-dependent; all pointers referencing the locations to be moved must also be changed. Clearly, this is not a trivial task and it is not clear that it can be done at all.

The result is that an operating system that restricts the placement of objects in a virtual address space will have a fragmented shared region that cannot be de-fragmented without enormous effort. Depending upon the degree of sharing this could mean a monotonically increasing shared region, which would be inimical to a 24x7 environment, i.e. one that is intended to be operative 24 hours a day, seven days a week.

3.3 PowerPC: Segmented translation

The IBM 801 introduced a segmented design that persisted through the POWER and PowerPC architectures [6, 29, 40, 54]; it is illustrated in Fig 1. Applications generate 32-bit “effective” addresses that are mapped onto a larger “virtual” address space at the granularity of *segments*, 256MB virtual regions. Sixteen segments comprise an application’s address space. The top four bits of the effective address select a segment identifier from a set of 16 registers. This segment ID is concatenated with the bottom 28 bits of the effective address to form an extended virtual address. This extended address is used in the TLB and page table. The operating system performs data movement and relocation at the granularity of pages, not segments.

The architecture does not use explicit address space identifiers; the segment registers ensure address space protection. If two processes duplicate an identifier in their segment registers they share that virtual segment by definition; similarly, protection is guaranteed if identifiers are *not* duplicated. If memory is shared through

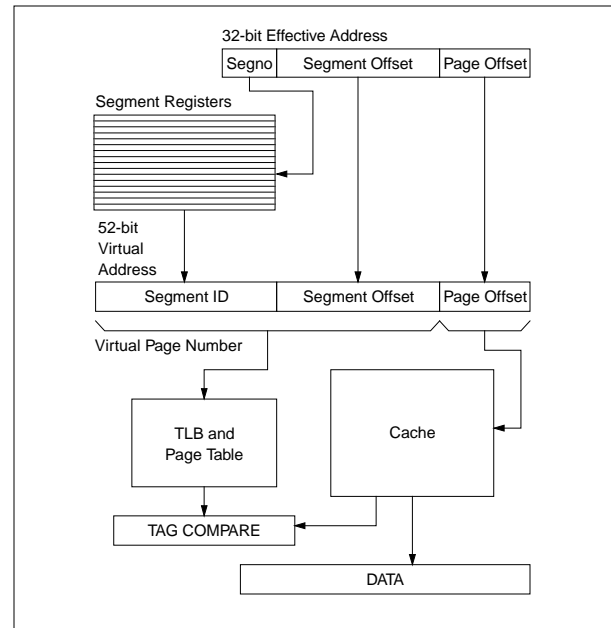


Figure 1: PowerPC segmented address translation. Processes generate 32-bit effective addresses that are mapped onto a 52-bit address space via sixteen segment registers, using the top four bits of the effective address as an index. It is this extended virtual address that is mapped by the TLB and page table. The segments provide address space protection and can be used for shared memory.

global addresses, no aliasing (and therefore no virtual-cache synonyms) can occur and the TLB and cache need not be flushed on context switch². This solution to the virtual cache synonym problem is similar to that of single address space operating systems—global addresses cause no synonym problems.

The relationship between the per-process effective address space and the global virtual address space is illustrated in Fig 2. This figure depicts how memory is shared in a segmented architecture using global addresses. This allows processes to map shared segments at arbitrary segment-aligned addresses, and to map shared segments at multiple locations if desired. However, since there is a one-to-one correspondence between physical addresses and global virtual addresses (the addresses used to reference the cache and TLB), there can be no synonym problems in a virtual cache; a physical datum can exist in one and only block of a virtual cache in any given instant, even if the cache is set associative.

2. Flushing is avoided until the system runs out of identifiers and must reuse them. For example, the address space identifiers on the MIPS R3000 and Alpha 21064 are six bits wide, with a maximum of 64 active processes [17, 33]. If more processes are desired, identifiers must be constantly reassigned, requiring TLB & virtual-cache flushes.

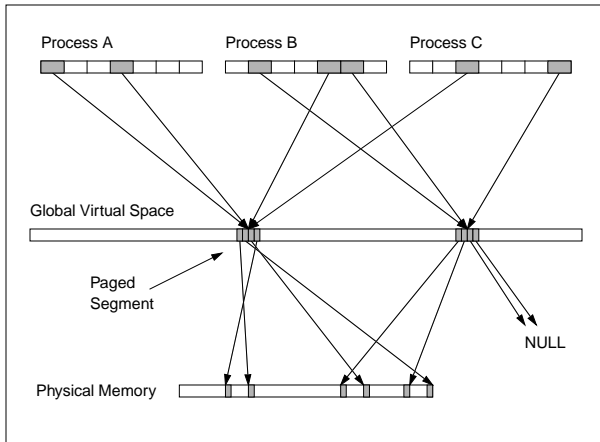


Figure 2: Virtual address aliasing in a segmented architecture. The figure shows three processes sharing two segments. None of the processes use the same virtual address for the same physical data, and two of the processes go so far as to map a segment at multiple locations within their address spaces. Nonetheless, these aliases will not result in any synonym problems in a virtual cache, since there is a one-to-one correspondence between pages in the global virtual space and pages in physical memory.

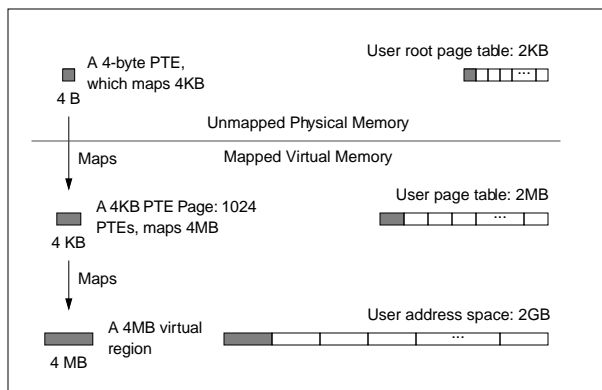


Figure 3: The MIPS 32-bit hierarchical page table. MIPS hardware provides support for a 2MB linear virtual page table that maps the 2GB user address space by constructing a virtual address from a faulting virtual address that indexes the mapping PTE in the user page table. This 2MB page table can easily be mapped by a 2KB user root page table.

3.4 MIPS: A simple 32-bit page table design

MIPS [23, 33] eliminated the page-table-walking hardware found in traditional memory management units, and in doing so demonstrated that software can table-walk with reasonable efficiency. It also presented a simple hierarchical page table design, shown in Fig 3. On a TLB miss, the hardware creates a virtual address for the mapping PTE in the user page table. The virtual page number (VPN) of the address that missed the TLB is used as an index into the user page table, which must be aligned on a 2MB virtual boundary. The base pointer, called *PTEBase*, is stored in a hardware register and is usually changed on context switch. This is illustrated as part of Fig 4. The advantage of this page table organization is that a small amount of wired-down memory

(2KB) can map an entire user address space efficiently; in the worst case, a user reference will require two additional memory lookups: one for the root-level PTE, one for the user-level PTE. The TLB miss handler is very efficient in the number of instructions it requires: the handler is less than ten instructions long, including the PTE load. We base our page table and cache miss design on this scheme for its simplicity and good performance.

3.5 SPUR: In-cache address translation

SPUR [26, 44, 56, 57] demonstrated that the storage slots of the TLB are not a necessary component in address translation. The architecture uses a virtually indexed, virtually tagged cache to delay the need for address translation until a cache miss occurs. On a miss, a hardware state machine generates the virtual address for the mapping PTE and searches the cache for that address. If this lookup misses, the state machine continues until the topmost level of the page table is reached, at which point the hardware requests the root PTE (at a known address) from physical memory.

The SPUR design eliminated specialized, dedicated hardware to store mapping information. However, it replaced the TLB with another specialized hardware translation mechanism—a finite state machine that searched for PTEs in general-purpose storage (the cache) instead of special-purpose storage (TLB slots).

3.6 VMP: Software-controlled caches

The VMP multiprocessor [11, 12, 13] places virtual caches under software control. Each processor node contains several hardware structures, including a central processing unit, a software-controlled virtual cache, a cache controller, and special memory. Objects the system cannot afford to have causing faults, such as root page tables and fault-handling code, are kept in a separate area called *local memory*, distinguished by the high-order bits of the virtual address. Code in local memory controls the caches; a cache miss invokes a fault handler that locates the requested data, possibly causes other caches on the bus to invalidate their copies, and loads the cache.

The scheme reduces the amount of specialized hardware in the system, including memory management unit and cache miss handler, and it simplifies the cache controller hardware. However, the design relies upon special memory that lies in a completely separate namespace from the rest of main memory.

4 Software-managed address translation

The *softvm* design requires a virtual cache hierarchy. There is no TLB, no translation hardware. When a refer-

ence fails to hit in the bottommost virtual cache a CACHEMISS exception is raised. We will refer to the address that fails to hit in the lowest-level cache as the *failing address*, and to the data it references as the *failing data*.

The general design is based on two observations. The first is that most high performance systems have reasonably large L2 caches, from 256KB found in many PCs to several megabytes found in workstations. Large caches have low miss rates; were these caches virtual, the systems could sustain long periods requiring no address translation at all. The second observation is that the minimum hardware necessary for efficient virtual memory is a software-managed cache miss at the lowest level of a virtual cache hierarchy. If software resolves cache misses, the operating system is free to implement whatever virtual-to-physical mapping it chooses. Wood demonstrated that with a reasonably large cache (128KB+) the elimination of a TLB is practical [56]. For the cache sizes we are considering, we reach the same conclusion (see the *Performance* section for details).

4.1 Handling the CACHEMISS exception

On a CACHEMISS exception, the miss handler loads the data at the failing address on behalf of another thread. The operating system must therefore be able to load a datum using one address and place it in the cache tagged with a different address. It must also be able to reference memory virtually or physically, cached or uncached; to avoid causing a cache-miss exception, the cache-miss handler must execute using physical addresses. These may be cacheable, provided that a cacheable-physical address that misses the cache causes no exception, and that a portion of the virtual space can be directly mapped onto physical memory.

When a virtual address misses the cache, the failing data, once loaded, must be placed in the cache at an index derived from the failing address and tagged with the failing address's virtual tag, otherwise the original thread will not be able to reference its own data. We define a two-part load, in which the operating system first specifies a virtual tag and set of protection bits to apply to the incoming data, then loads the data with a physical address. The incoming data is inserted into the caches with the specified tag and protection information. This scheme requires two privileged instructions to be added to the instruction set architecture (ISA)³: SPECIFYVTAG and LOAD&MAP, depicted in Fig 4.

3. Many ISAs leave room for such management instructions, e.g. the PowerPC ISA *mtspr* and *mfspr* instructions (move to/from special purpose register) would allow implementations of both functions.

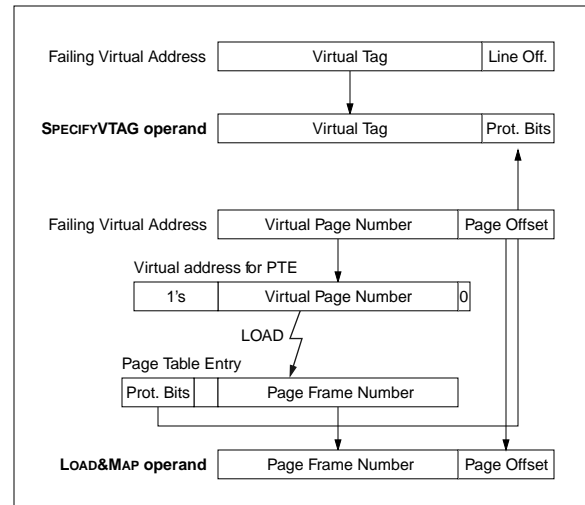


Figure 4: SPECIFYVTAG and LOAD&MAP. The top figure illustrates SPECIFYVTAG, the bottom figure illustrates LOAD&MAP. The LOAD&MAP example resembles a MIPS-style page table lookup; the PUMA global page table is very similar to the MIPS organization. However, while the user-level page table in MIPS is located at *PTEBase*, a variable stored in a hardware register, the PUMA global page table is located at a known offset at the top of the global address space, thus the 1's in the most significant bits of the virtual address for the PTE.

SPECIFYVTAG instructs the cache to insert future incoming data at a specific offset in the cache, tagged with a specific label. Its operand has two parts: the virtual tag (VTAG) comes from the failing virtual address; the protection bits come from the mapping PTE. The bottom half of the VTAG identifies a block within the cache, the top half is the tag. Note that the VTAG is larger than the virtual page number; the hardware should not assume *any* overlap between virtual and physical addresses beyond the cache line offset. This is essential to allow a software-defined page size.

The operand of a LOAD&MAP is a physical or virtual address. The datum identified by the operand is loaded from the cache or memory and then (re-) inserted into the cache at the cache block determined by the previously executed SPECIFYVTAG, and tagged with the specified virtual tag. Thus an operating system can translate data that misses the cache, load it from memory (or even another location in the cache), and place it in any cache block, tagged with any value. When the original thread is restarted, its data is in the cache at the correct line, with the correct tag. Note the operations can be performed out of order for performance reasons, as long as the tag arrives at the cache/s before the data arrives. Note also that without hardware support, the two-part load must not be interrupted by another two-part load.

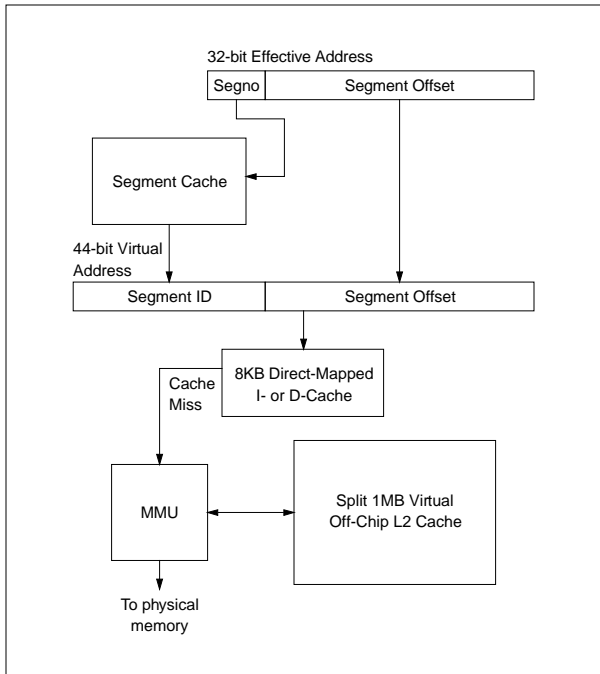


Figure 5: The PUMA address translation mechanism. The processor uses PowerPC-style segmentation to translate a user 32-bit effective address into an extended 44-bit virtual address, used to reference both caches in the virtual hierarchy. However, unlike the 801 and PowerPC, more than 4 bits are used to index the segment cache; the top 10 bits indicate a segment identifier, dividing the 4GB user address space into 1024 4MB virtual segments.

5 The PUMA memory management architecture

5.1 Hardware

The PUMA is a high clock-rate 32-bit PowerPC. Its memory management architecture is based on PowerPC segmentation and provides a software-managed address translation mechanism [32]. Processes generate 32-bit addresses which are extended through the segmentation mechanism to 44-bit addresses. The entire cache hierarchy is virtual, and when a reference misses in the L2 cache, the miss is handled in software, not hardware. A TLB is not necessary, as protection information is kept in the caches. As we have shown [32], the minimum hardware necessary for efficient virtual memory is a virtual cache hierarchy with a software-managed L2 cache miss, and many other studies have already shown that a TLB is not necessary for good performance [14, 44, 56, 57]. The benefit of a software-managed cache miss is that the operating system is free to implement whatever mapping organization it chooses.

The PUMA address translation mechanism is shown in Fig 5. Both caches in the two-level virtual cache hierarchy are split and virtual. The L1 cache is 8KB/8KB, the L2 cache is 512KB/512KB. When references miss in the L2 caches the MMU generates a `CACHEMISS` excep-

tion; the operating system is responsible for handling address translation and loading the data at the faulting address. The operating system can reference memory virtually or physically, depending on the segment identifier used. As in MIPS [33], virtual memory is divided into regions of different behavior; two of the regions map directly onto physical memory, the rest are virtual. The MMU checks the top twelve bits of the 44-bit virtual address; the patterns 0x000 and 0x001 represent cached and uncached physical memory, respectively. This places a limit of 2^{32} bytes on the physical memory size—exactly the amount of space that the operating system can address at once. All other segment identifiers represent virtual segments and can fault. When the MMU receives virtual addresses prefixed with either of the reserved patterns, it treats them specially. If the pattern is 0x001, the MMU first checks the L2 cache. If the data is not found in the cache, or if the pattern is 0x000, the MMU passes the rest of the address on to physical memory. When the data returns from main memory it is sent to the FXU; if the pattern is 0x001, the MMU also inserts the data into the cache.

To handle a `CACHEMISS` exception, the operating system must load the data at the faulting address on behalf of the user process. This is complicated by the fact that the data must be accessed using its physical name but stored in the caches by its virtual name. The architecture defines a two-part load, in which the operating system first sends a virtual page number to the MMU, then loads the data with the physical address, instructing the MMU to store it in the cache tagged with the previously indicated VPN. Thus when the faulting thread is restarted, its data is in the cache at the correct line, with the correct tag. These functions are provided by two new privileged instructions: `SPECIFYVTAG` and `LOAD&MAP`.

Note that the hardware does not manage the 0x000/0x001 overlap problem. For instance, if one writes one datum to location 0x001000BABA0 and then a different datum to location 0x000000BABA0 (the same location, just not cached), the original datum will still be left in the cache. This will cause problems when a subsequent cacheable read or the inevitable writeback occurs. It is up to the operating system to avoid this behavior, which can be done by careful organization of cached/uncached regions.

Segment protection information is kept in the segment registers, and page protection information is kept in the global page table. An entire object is mapped into an address space with one protection, but it can have different segment-level protections every time it is mapped (either in another address space or at another offset in the same address space). The global page table keeps page-protection information, which is maintained

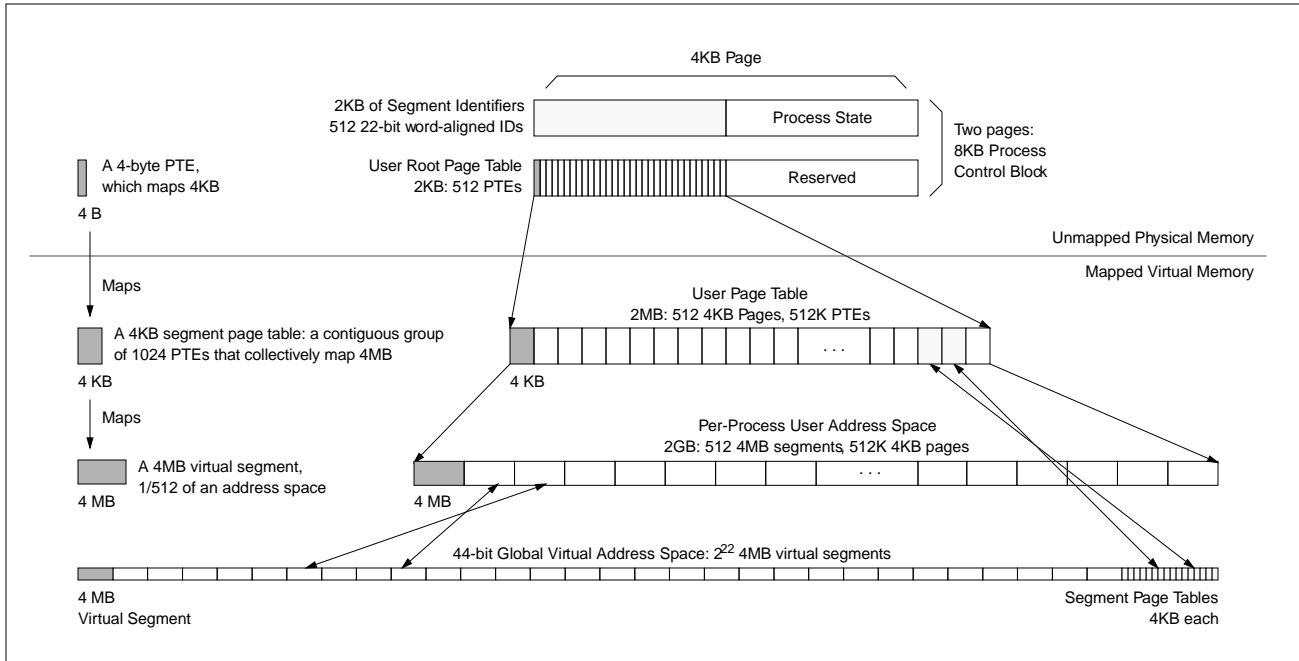


Figure 6: The disjunct page table organization. There is a single linear page table at the top of the 44-bit address space that maps the entire address space. The 4KB PTE pages (segment page tables) of the user page table are taken directly from this global page table. Therefore, though it may seem that there is a separate user page table for every process, each page table is simply mapped onto the global space; the only per-process allocation is for the user root page table. Though it is drawn as an array of contiguous pages, the user page table is really a disjunct set of 4KB pages in the global space.

in hardware on a per-cacheline basis. Keeping the page-protection bits consistent across multiple cachelines is not trivial; this is discussed in more detail in the *Performance* section.

5.2 Software

It is not yet clear how finely-segmented the PUMA architecture will be; the granularity depends upon resource constraints, both in the process and the design team. For the purposes of this report, we assume the best case and describe the design for a PUMA processor that divides the user address space into 1024 4MB segments.

We use a global page table to guarantee one-to-one mappings between virtual space and physical space. We use the PUMA segment cache to allow processes to map segments at arbitrary locations. The per-process address space is 2GB; like MIPS [33], the top half of the address space is owned by the kernel—although here it is a software convention, not a hardware stipulation.

The global page table is a linear structure at the top of the 2^{44} -byte global address space. The page table is 2^{34} bytes long (pages are software-defined at 4KB, PTEs are 4B). The table, shown in Fig 6, has a two-tiered hierarchical organization, accessed bottom-up. The top tier is a 2KB structure wired down in physical memory while the process is running. This per-process root page table occupies one half of a page—part of the process control block. The lower tier of the page table is a 2MB

structure in virtual space. It is divided into 512 4KB regions, each of which (collectively) maps one of the 4MB segments of the user’s address space. These PTE pages are called *segment page tables* since each maps a virtual segment. They are located in the topmost portion of the global virtual space.

The page fault algorithm is shown in Fig 7. Processes generate 32-bit effective addresses that are extended to 44 bits by segmentation, replacing the top four bits of the effective address. In step 1 of the figure, the operating system uses the faulting virtual address to construct the virtual address of the page table entry that maps the faulting address. The VPN of a 44-bit failing global virtual address is used as an index into the global page table to reference the PTE mapping the failing data (the UPTE). The bottom two bits of the address are 0’s, since the PTE size is four bytes. The top ten bits of the address are 1’s since the table is at the very top of the global space. This virtual address is sent to the L1 cache, and the L2 cache if there is an L1 miss.

If this misses in the L2 cache, the operating system takes a recursive *CACHEMISS* exception. At this point, we must locate the mapping PTE in the user root page table. This table is an array of PTEs that cannot be indexed by a global VPN. It mirrors the structure of the user’s perceived address space, not the structure of the global address space. Therefore it is indexed by a portion of the original 32-bit effective address. The top 10 bits of the effective address index 1024 PTEs that would

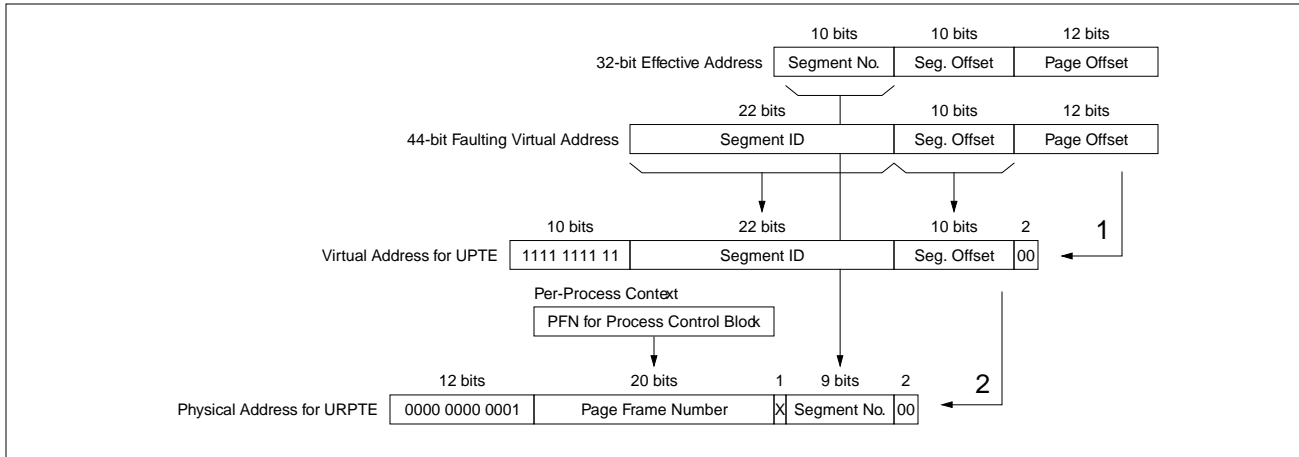


Figure 7: The PUMA page fault algorithm. Step 1 is the result of a user-level L2 cache miss; the operating system builds a virtual address for a PTE in the global page table. If this PTE is not found in the L1 or L2 cache a root PTE is loaded, shown in step 2. One special requirement is a register holding the initial faulting address. Another required hardware structure, the per-process context register, points to the process control block of the active process.

map a 4MB user page table, which would in turn map a 4GB address space. Since the top bit of the effective address is guaranteed to be zero (the address is a user reference), only the bottom nine bits of the top ten are meaningful; these bits index the array of 512 PTEs in the user root page table. In step 2, the operating system builds a physical address for the appropriate PTE in the user root page table (the URPTE), a 44-bit virtual address whose top 12 bits are 0x001, indicating that it is physical and cacheable. Note that this 44-bit address can be used to index the virtual caches, even though it is equivalently mapped to a physical address. The operating system then loads the URPTE, which maps the UPTE that missed the cache at the end of step 1. When control is returned to the miss handler in step 1, the UPTE load retry will complete successfully.

Once an appropriate mapping structure has been loaded, the operating system needs to load the data at the faulting user address and place it in the virtual cache. It uses the two-part load procedure described earlier. First an appropriate VPN is constructed from the 44-bit faulting virtual address and handed to the MMU as an operand of a SPECIFYVTAG instruction. The operating system then performs a LOAD&MAP using the physical address for the failing data, built from the PFN in the UPTE and the page offset from the failing address. The MMU places the data in the cache at an offset determined by the specified VTAG. This loads the failing data and inserts it into the cache using the user's virtual tag. When the faulting instruction is restarted, the data is in the cache hierarchy at the correct offset, tagged with the correct VPN.

Note that the 'X' bit between the page frame number and the segment number in the physical address is reserved. To conserve space, we double up on user root

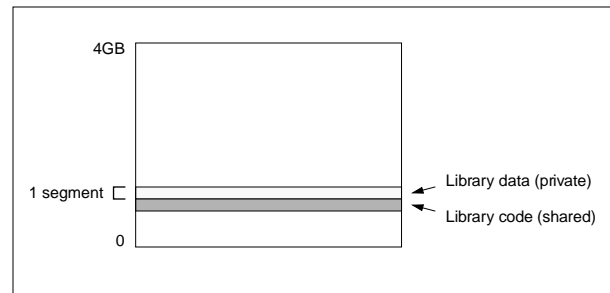


Figure 8: The location of shared libraries in process address spaces. The code and data of a library are separated into two halves, each occupying its own virtual segment. The code segment can be shared by multiple processes, but the data segment must be private for every process. The code and data segments are located at the same virtual address in all processes, so the code can use absolute addresses to reference locations in the data segment.

page tables, and put two URPTs into a single page. This bit may be zero or one.

5.3 Dynamically linked shared libraries

To give a feeling for how we use the segmented space to implement shared memory, we present a brief overview of our implementation of dynamically linked shared libraries. We do not require position independent code, as we amortize one-time library linking across multiple program uses. We do not statically fix library addresses, but like OS/2 we require that libraries be loaded into every address space at the same virtual location. Therefore, if a process does not use a particular library the space is unusable.

Each library has code and data segments; code segments are shared, data segments are private. We do not allow shared data segments, otherwise libraries would have to be reentrant. The simple shared library organization is shown in Fig 8. Libraries occupy two or more (possibly adjacent) segments in the user's address

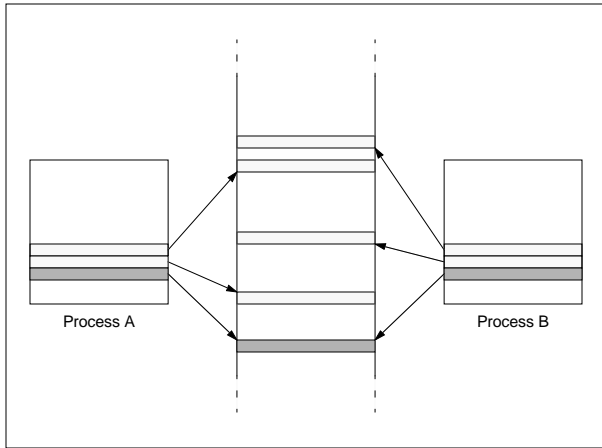


Figure 9: Two processes sharing a library. The code segment is shared between the two processes but the data segments are private. The data segments are located at the same virtual address in both processes, so the code can use the same virtual address to reference the appropriate data. The segmentation mechanism provides the ability to use different aliases for different physical data.

space. The code segment is shared by all processes that use the library. The data segment is mapped into every process at the same location, but every process uses a different segment identifier to keep from accessing each other's data. Library code can therefore use absolute pointers to reference internal data and code.

An example of two applications sharing a library is shown in Fig 9. The figure illustrates an example of a library with large data requirements; the data region cannot fit into a single hardware segment, so it stretches across two segments. As the figure shows, segments that are next to each other in the per-process address spaces may end up nowhere near each other in the global space. Despite this, libraries and shared regions in general can store and share pointers in the global space.

The advantage of this organization is that library code can use absolute values for both jump offsets and data loads and stores. The code will be shared and the data private, and both will be at the same offsets in all processes. This is a simple implementation of position independent code.

The disadvantage is that the library code will need to be linked at the time of execution, but only for the first process to use the library. Load/store addresses and jump offsets will need to reflect what segment numbers the library code and data have been loaded into. All libraries are created by the compiler based at address zero, and the code or data segment offset can simply be added to all absolute addresses. To speed up system response time, linking can be performed on a page-by-page basis when the page is first demand-faulted in.

Alternatively, the operating system can reserve segment slots for certain libraries, as seen in several com-

mercial operating systems, including Solaris, SVR3, and VMS.

5.4 Memory system requirements, revisited

We now revisit the memory management requirements listed earlier, and discuss how the PUMA design supports them.

Address space protection and large address spaces.

These are satisfied in our example through the use of PowerPC segments in conjunction with a virtual cache. As described earlier, segments provide address space protection, and by their definition provide a global virtual space onto which all effective addresses are mapped. A process could use its 4GB space as a window onto the larger space, moving virtual segments in and out of its working set as necessary. This type of windowing mechanism is used on the PA-RISC [27].

Shared memory. The sharing mechanism is defined by the page table. Sharing memory via global addresses can simplify virtual cache management; this is the scheme used in many systems [7, 8, 16, 18, 20, 21, 46], and has been shown to have good performance.

Fine-grained protection. One can maintain protection bits in the cache, or in an associated structure like a TLB. If one could live with protection on a per-segment basis, one could maintain protection bits in the segment registers. We maintain protection bits in the cache line. Protection granularity therefore becomes a software issue; the page size can be anything from the entire address space down to a single cache line. Note the choice of this granularity does not preclude us from implementing segment-level protection as well. The disadvantage is that if one chooses a page size larger than a single cache line, protection information must be replicated across multiple cache lines and the operating system must manage its consistency. We analyze this later.

Sparse address spaces. Sparse address space support is largely a page table issue. Hardware can either get out of the way of the operating system and allow any type of page table organization, or it can inhibit support for sparse address spaces by defining a page table organization that is not necessarily suitable. By eliminating translation hardware, one frees the operating system to choose the most appropriate structure.

Table 1: Qualitative comparison of cache-access/address-translation mechanisms

Event	Frequency of Occurrence		Actions Performed by Hardware and Operating System per Occurrence of Event	
	I-side	D-side	TLB + Virtual cache	Software-Mgd Addr Translation
L1 hit, TLB hit	96.7%	95.8%	L1 access (w/ TLB access in parallel)	L1 access
L1 hit, TLB miss	0.01%	0.06%	L1 access + page table access + TLB reload	L1 access
L1 miss, L2 hit, TLB hit	3.2%	3.9%	L1 access + L2 access	L1 access + L2 access
L1 miss, L2 hit, TLB miss	0.03%	0.09%	L1 access + page table access + TLB reload + L2 access	L1 access + L2 access
L1 miss, L2 miss, TLB hit	0.008%	0.12%	L1 access + L2 access + memory access	L1 access + L2 access + page table access + memory access
L1 miss, L2 miss, TLB miss	0.0001%	0.0009%	L1 access + page table access + TLB reload + L2 access + memory access	L1 access + L2 access + page table access + memory access

Superpages. By removing the TLB one removes hardware support for superpages, but as with sparse address spaces one also frees the operating system to provide support through the page table. For instance, a top-down hierarchical page table (as in the x86 [31]) would provide easy support for superpages. A guarded page table [38, 39] would also provide support, and would map a large address space more efficiently, as would the inverted page table variant described by Talluri, et al. [49].

Direct memory access. While our design provides no explicit support for DMA, and actually makes DMA more difficult by requiring a virtual cache, direct memory access is still possible. We perform DMA by flushing affected pages from the cache before beginning a transfer, and restricting access to the pages during transfer.

6 Performance of the design

Many studies have shown that significant overhead is spent servicing TLB misses [1, 4, 9, 28, 42, 45, 48]. In particular, Anderson, et al. [1] show TLB miss handlers to be among the most commonly executed primitives, Huck and Hays [28] show that TLB miss handling can account for more than 40% of total run time, and Rosenblum, et al. [45] show that TLB miss handling can account for more than 80% of the kernel’s computation time. Typical measurements put TLB handling at 5-10% of a normal system’s run time.

The question that we have asked is *does the TLB buy us anything?* Do its benefits outweigh its overhead? We discovered that with the large caches available today (most of which are larger than the entire memory systems of computers twenty to thirty years ago when virtual memory was invented and when TLBs came into fashion), address translation can be performed much more efficiently *without* a TLB. It does not, in fact, buy us anything.

6.1 A qualitative overview

The SPUR and VMP projects demonstrated that with large virtual caches the TLB can be eliminated with no performance loss, and in most cases a performance gain. For a qualitative, first-order performance comparison, we enumerate the scenarios that a memory management system would encounter. These are shown in Table 1, with frequencies obtained from SPECint95 traces on a PowerPC-based AIX machine (frequencies do not sum to 1 due to rounding). The model simulated has 8K/8K direct-mapped virtual L1 caches (in the middle of the L1 cache sizes simulated), 512K/512K direct-mapped virtual L2 caches (the smaller of the two L2 cache sizes simulated), and a 16-byte linesize in all caches. As later graphs will show, the small linesize gives the worst-case performance for the software-managed scheme. The model includes a simulated MIPS-style TLB [33] with 64 entries, a random replacement policy, and 8 slots reserved for root PTEs.

The table shows what steps the operating system and hardware take when cache and TLB misses occur. Note that there is a small but non-zero chance a refer-

ence will hit in a virtual cache but miss in the TLB. If so, the system must take an exception and execute the TLB miss handler before continuing with the cache lookup, despite the fact that the data is in the cache. On TLB misses, a software-managed scheme should perform much better than a TLB scheme. When the TLB hits, the two schemes should perform similarly, except when the reference misses the L2 cache. Here the TLB already has the translation, but the software-managed scheme must access the page table for the mapping (note that the page table entry may in fact be cached). Software-managed translation is not penalized by placing PTEs in the cache hierarchy; many operating systems locate their page tables in cached memory for performance reasons.

6.2 Baseline overhead

Table 2 shows the overheads of TLB handling in several operating systems as percent of run-time and CPI. Percent of run-time is the total amount of time spent in TLB handlers divided by the total run-time of the benchmarks. CPI overhead is the total number of cycles spent in TLB handling routines divided by the total number of cycles in the benchmarks. The data is taken from previous TLB studies [4, 41, 42] performed on MIPS-based DECstations, which use a software-managed TLB. CPI is not directly proportional to run-time overhead for two reasons: (1) the run-time overhead contains page protection modifications and the CPI overhead does not, and (2) memory stalls make it difficult to predict total cycles from instruction counts.

Table 2: TLB overhead of several operating systems

System	Overhead (% run-time)	Overhead (CPI)
Ultrix	2.03%	0.042
OSF/1	5.81%	0.101
Mach3	8.21%	0.162
Mach3+AFSin	7.77%	0.220
Mach3+AFSout	8.88%	0.281

Table 3 gives the overheads of the software-managed design, divided by benchmark to show a distribution. The values come from trace-driven simulation of the SPEC95 integer suite. The simulations use the same middle-of-the-line cache organization as before (8K/8K L1, 512K/512K L2, 16-byte linesize throughout), but replace the TLB with software address translation. Our memory penalties are 1 cycle to access the L1 cache, 20 cycles to access the L2 cache, and 90 cycles to access main memory.

Table 4 gives a more detailed breakdown of costs for one of the benchmarks: gcc. Our example miss handler from the previous section requires 10 instructions including two loads. It is very similar to the MIPS TLB

Table 3: Overhead of software-managed address translation

Workload	Overhead (CPI)
m88ksim	0.003
li	0.003
go	0.004
compress95	0.009
perl	0.019
ijpeg	0.052
vortex	0.060
gcc	0.097
Weighted Average:	0.033

refill handler that requires less than 10 instructions including one load, taking 10 cycles when the load hits in the cache, or 40+ when the load misses in the cache, thereby forcing the reference to main memory [4]. In our model, the L2 cache miss handler always takes 10 cycles, and runs whenever we take an L2 cache miss (labeled *L2 I-Cache miss* or *L2 D-Cache miss* in the table). When the PTE load in the handler misses the L1 cache (*Miss handler L1 D-miss*) we take an additional 20 cycles to go to the L2 cache to look for the PTE. If that load misses we either take a recursive cache miss (if handling a user-miss, therefore the PTE address is virtual, accounted for in *L2 D-Cache miss*), or the address is physical and goes straight through to main memory (*Miss handler L2 D-miss*, 90 cycles). When the miss handler is handling a miss from the handler itself, we need also load the failing UPTE on behalf of the handler (*Miss handler Load UPTE*, 90 cycles).

Table 4: Breakdown of GCC overhead

Event	Frequency (per instr.)	Penalty per Occurrence	Overhead (CPI)
L2 D-Cache Ld/St miss	0.000697	10 cycles	0.006970
L2 I-Cache I-fetch miss	0.004756	10 cycles	0.047560
Miss handler L1 D-miss	0.000596	20 cycles	0.011920
Miss handler L2 D-miss	0.000032	90 cycles	0.002880
Miss handler Load UPTE	0.000053	90 cycles	0.004770
Miss handler L1 I-miss	0.000985	20 cycles	0.019700
Miss handler L2 I-miss	0.000035	90 cycles	0.003150
Total CPI:			0.096950

Additionally, the miss-handler code can miss in the L1 or L2 I-caches; since it is mapped directly onto physical memory it does not cause a cache miss itself. How-

ever, for every instruction fetch that misses in the L1 cache we take a 20-cycle penalty to reference the L2 cache; for every L2 miss we take a 90-cycle penalty to reference physical memory.

The average overhead of the scheme is 0.033 CPI. This is about the overhead measured of Ultrix on MIPS, considered to be an example of an efficient match between OS and architecture. This CPI is several times better than that of Mach, which should result in a runtime savings of at least 5% over Mach.

6.3 Writebacks

We keep the translation of a virtual address in the cache line with the data. The translation is only needed in the L2 cache, since we maintain inclusion between the two caches. This simplifies writeback, as it allows the L2 cache to perform writebacks without the aid of the CPU, without interrupting the processor. However, the operating system must ensure that the translations kept in the cache never become stale; whenever a virtual page is re-mapped, the operating system must modify or flush any entries belonging to the affected page.

6.4 Fine-grained protection

As mentioned earlier, managing protection information can be inefficient if we store protection bits with each cache line. If the protection granularity is larger than a cache line, the bits must be replicated across multiple lines. Keeping the protection bits consistent across the cache lines can cause significant overhead if page protection is modified frequently. The advantage of this scheme is that the choice of protection granularity is completely up to the operating system. In this section, we determine the overhead.

Table 5: Page protection modification frequencies in Mach3

Workload	Page Protection Modifications	Modifications per Million Instructions
compress	3635	2.8
jpeg_play	12083	3.4
IOzone	3904	5.1
mab	27314	15.7
mpeg_play	26129	19.0
gcc	35063	22.3
ousterhout	15361	23.8
	Weighted Average:	11.3

We performed a study on the frequency of page protection modifications in the Mach operating system. The benchmarks are the same as in [42], and the operating system is Mach3. We chose Mach as it uses copy-on-write liberally, producing 1000 times the page-protec-

tion modifications seen in Ultrix [42]. We use these numbers to determine the protection overhead of our system; this should give a conservative estimate for the upper bound. The results are shown in Table 5.

Page-protection modifications occur on the average of 11.3 for every million instructions. At the very worst, for each modification we must sweep through a page-sized portion of the L1 and L2 caches to see if lines from the affected page are present. Overhead therefore increases with larger page sizes (a software-defined parameter) and with smaller linesizes (a hardware-defined parameter). On a system with 4KB pages and a 16-byte linesize, we must check 256 cache lines per modification. Assuming an average of 10 L1 cache lines and 50 L2 caches lines affected per modification⁴, if L1 cache lines can be checked in 3 cycles and updated in 5 cycles (an update is a check-and-modify), and L2 cache lines can be checked in 20 cycles and updated in 40 cycles, we calculate the overhead as follows. Of 256 L1 cache lines, 10 must be updated (5 cycles), the remaining 246 need only be checked (3 cycles); of 256 L2 cache lines 50 must be updated (40 cycles), the remaining 206 need only be checked (20 cycles); the overhead is therefore 6908 cycles per page-protection modification ($10 * 5 + 246 * 3 + 50 * 40 + 206 * 20$). This yields between 0.019 and 0.164 CPI ($6908 * 2.8 * 10^{-6}$ and $6908 * 23.8 * 10^{-6}$). This is in the range of Ultrix and OSF/1 overheads and at the lower end of Mach’s overhead. This translates to a worst case of 2-7% total execution time. If the operating system uses page-protection modification as infrequently as in Ultrix, this overhead decreases by three orders of magnitude to 0.0001 CPI, or about 0.01% execution time.

We can improve this by noting that most of these modifications happen during copy-on-write. Often the protections are being increased and not decreased, allowing one to update protection bits in each affected cache line lazily—to delay an update until a read-only cache line is actually written, at which point it would be updated anyway.

6.5 Sensitivity to cache organization

The graphs in Fig 10 show the sensitivity of software-managed address translation to cache size and cache linesize. The benchmarks are from SPEC95 as before; we only show graphs for the two worst-performing benchmarks—gcc and vortex. The numbers differ slightly

4. We chose these numbers after inspecting individual SPEC95 benchmark traces, which should give conservative estimates: (1) SPEC working sets tend to be smaller than normal programs, resulting in less page overlap in the caches, and (2) individual traces would have much less overlap in the caches than multiprogramming traces.

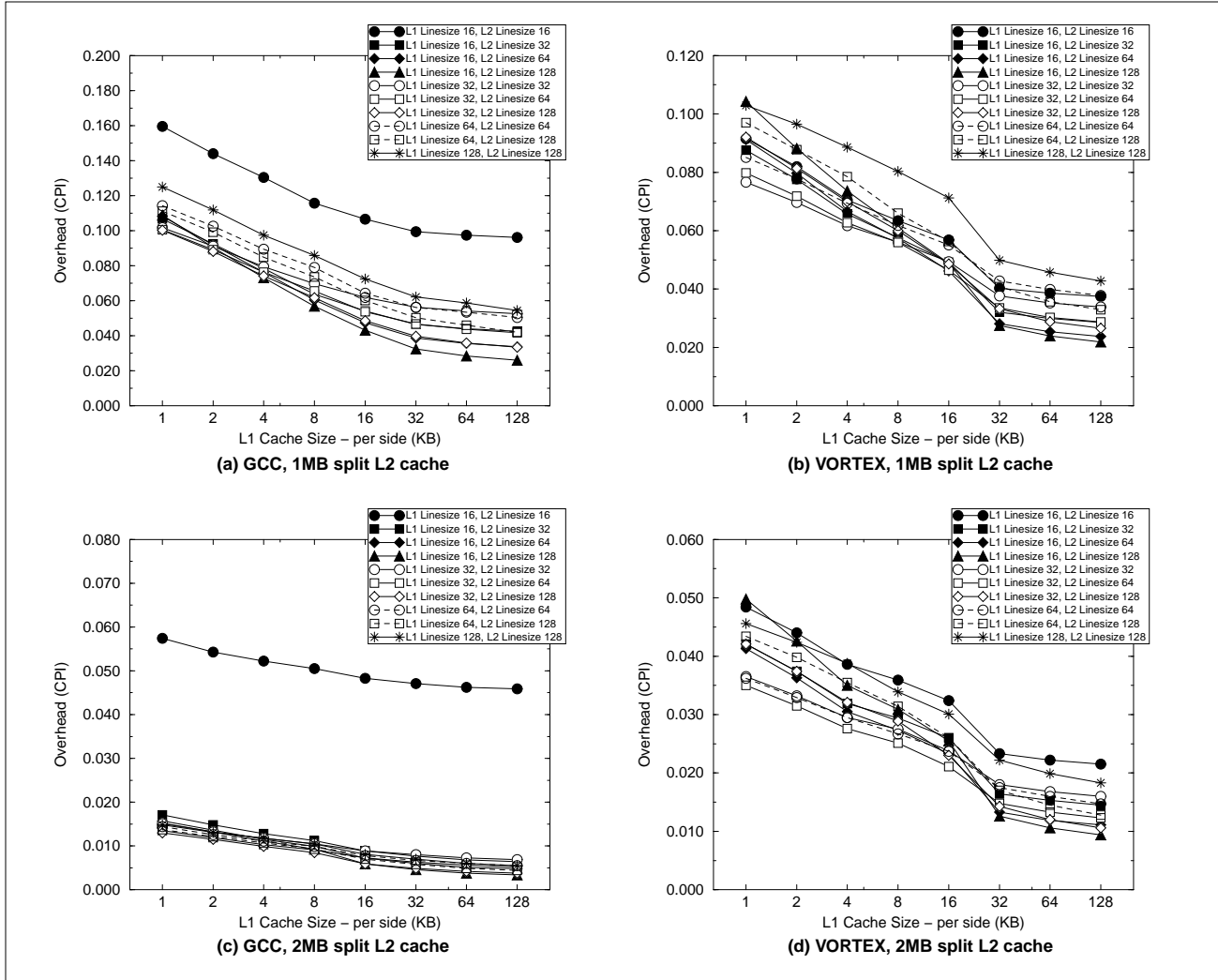


Figure 10: The effect of cache size and linesize on software-managed address translation. The figure shows two benchmarks—gcc and vortex. All caches are split. L1 cache size is varied from 1K to 128KB per side (2KB to 256KB total), and L2 cache size is varied from 512KB to 1024KB per side (1MB to 2MB total). Linesizes are varied from 16 bytes to 128 bytes; the L2 linesize is never less than the L1 linesize. In each simulation, the I-caches and D-caches have identical configurations. We apologize for using different y-axis scales; however, they better show the effects of linesize for a given cache size.

from those presented in Table 3; the benchmarks were not run to completion for this study, but were stopped after 1 billion references for each.

Besides the familiar signature of diminishing returns from increasing linesize (e.g., the two largest overheads in Fig 10a are from the smallest and largest linesizes), the graphs show that cache size has a significant impact on the overhead of the system. For gcc, overhead decreases by an order of magnitude when the L2 cache is doubled, and decreases by a factor of three as the L1 cache increases from 1KB to 128KB (2KB to 256KB total L1 cache size); for vortex, overhead decreases by a factor of two as the L2 cache doubles, and decreases by a factor of three as L1 increases from 1KB to 128KB. Within a given cache size, linesize choice can affect performance by a factor of two or more (up to ten for some configurations).

The best organization should result in an overhead an order of magnitude lower than that calculated earlier—to less than 0.01 CPI, or a run-time overhead far less than 1%. This suggests that software-managed address translation is viable today as a strategy for faster, nimbler systems.

7 Page table efficiency

The theoretical minimum page table size is 0.1% of working set size, assuming 4KB pages, 4B page table entries, and fully-populated page table pages. However, most virtual memory organizations do not share PTEs when pages are shared; in many operating systems, for every shared page there is more than one PTE in the page tables. Khalidi & Talluri show that these extra

PTEs can increase the page table size by an order of magnitude or more [34].

We compare the size of the global page table to the theoretical minimum size of a traditional page table. Khalidi & Talluri report that the average number of mappings per page on an idle system is 2.3, and the average number of mappings to *shared* pages is 27. This implies that the ratio of private to shared pages in an average system is 19:1 or that 5% of a typical system's pages are shared pages.⁵ These figures are used in our calculations. The overhead of a traditional page table (one in which there must be multiple PTEs for multiple mappings to the same page) can be calculated as:

$$\frac{(\text{number of PTEs})(\text{size of PTE})}{(\text{number of pages})(\text{size of page})} = \frac{(p + 27s)4}{(p + s)4096} = \frac{(p + 27s)}{(p + s)1024}$$

where p is the number of private (non-shared) pages in the system, and s is the number of shared pages in the system. We assume a ratio of 1024:1 between page size and PTE size. This represents the theoretical minimum overhead since it does not take into account partially-filled PTE pages. For every shared page there is on average 27 processes mapping it, therefore the page table requires 27 PTEs for every shared page. The overhead is in terms of the physical-page working set size; the fraction of physical memory required to map a certain number of physical pages. As the percentage of sharing increases, the number of physical pages does not increase, but the number of PTEs in the page table does increase.

The global page table overhead is calculated the same way, except that PTEs are not duplicated when pages are shared. Thus, the overhead of the table becomes a constant:

$$\frac{(p + s)4}{(p + s)4096} = \frac{1}{1024}$$

Clearly, the global page table is smaller than a traditional page table, and approaches the minimum size necessary to map a given amount of physical memory. Fig 11 shows the overhead of each page table organization as the level of sharing in a system changes. In an average system, where 5% of the pages are shared, we should expect to use less than half the space required by a traditional page table organization.

5. The average number of mappings per page is the total number of mappings in the system divided by the total number of pages, or $\frac{p + 27s}{p + s} = 2.3$, yielding a $p:s$ ratio of 19:1.

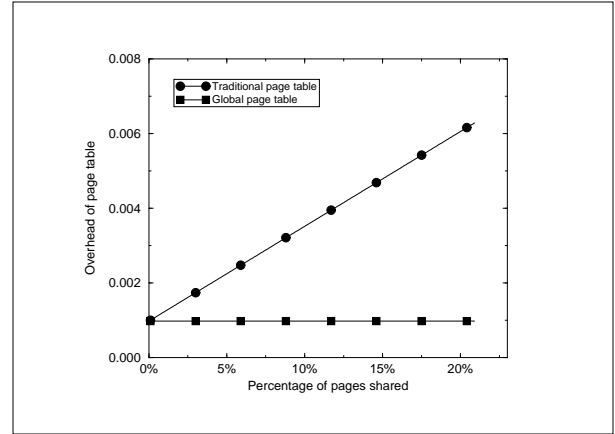


Figure 11: Comparison of page table space requirements. The x-axis represents the degree of sharing in a system, as the number of pages that are shared ($s/(p + s)$). The y-axis represents the overhead of the page table, as the size of the page table divided by the total size of the data pages. Khalidi and Talluri's research shows that in an average system, 5% of the pages are shared. Therefore, we would expect a disjunct page table to require slightly less than half the space needed by a more traditional page table.

8 Summary

We are building a high clock-rate 32-bit PowerPC. In order to meet the memory requirements of a high-speed processor and avoid the potential slowdown of address translation, we employ a two-level virtual cache hierarchy. Virtual caches help achieve fast clock speeds but have traditionally been left out of microprocessor architectures because they have the potential for data inconsistencies, requiring significant management overhead. A segmented architecture adds a second level of indirection and allows a system to use a virtual cache organization without explicit consistency management, as long as the operating system ensures that there is a one-to-one mapping of pages between the segmented address space and physical memory. We use a *disjunct page table* organization to maintain a one-to-one mapping while allowing processes to map virtual segments at different locations in their address spaces, or even at multiple locations in their address spaces. The disjunct page table organization requires less than half the space of a more traditional page table organization.

For the design of the memory management system, we have returned to first principles and discovered a small set of hardware structures that provide support for address space protection, shared memory, large sparse address spaces, and fine-grained protection at the cache-line level. This set does not include address-translation hardware; we show that address translation can be managed in software efficiently. Current virtual memory systems such as Mach exact an overhead of 0.16 to 0.28 cycles per instruction to provide address

translation; our software scheme requires 0.03 CPI (2% run-time, with a 16KB L1 cache, and 1MB L2), about the same as the overhead of Ultrix on MIPS. If copy-on-write and other page-protection modifications are used as frequently as in Mach, protection-bit management can increase this overhead to that of OSF/1 or Mach. However, the number of page-protection modifications in Ultrix represent a negligible overhead. With slightly larger caches (2MB L2, common in today's systems), the overhead of software-managed address translation should reduce to far less than 1% of run-time. Therefore software-managed address translation is a viable strategy for high-end computing today, achieving better performance with less hardware.

Beyond the performance gains suggested by our simulations, the benefits of a minimal hardware design are three-fold. First, moving address translation into software creates a simpler and more flexible interface; as such it supports much more innovation in the operating system than would a fixed design. Second, a reduction in hardware will leave room for more cache structures, increasing performance. Last, simpler hardware should be easier to design and debug, cutting down on development time.

References

- [1] T. E. Anderson, H. M. Levy, B. N. Bershad, and E. D. Lazowska. "The interaction of architecture and operating system design." In *Proc. Fourth Int'l Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS-4)*, April 1991, pp. 108–120.
- [2] A. W. Appel and K. Li. "Virtual memory primitives for user programs." In *Proc. Fourth Int'l Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS-4)*, April 1991, pp. 96–107.
- [3] M. J. Bach. *The Design of the UNIX Operating System*. Prentice-Hall, Inc., Englewood Cliffs, NJ, 1986.
- [4] K. Bala, M. F. Kaashoek, and W. E. Wehl. "Software prefetching and caching for translation lookaside buffers." In *Proc. First USENIX Symposium on Operating Systems Design and Implementation*, November 1994.
- [5] B. N. Bershad, C. Chambers, S. Eggers, C. Maeda, D. McNamee, P. Pardyak, S. Savage, and E. G. Sizer. "SPIN – an extensible microkernel for application-specific operating system services." Tech. Rep. 94-03-03, University of Washington, February 1994.
- [6] A. Chang and M. F. Mergen. "801 storage: Architecture and programming." *ACM Transactions on Computer Systems*, vol. 6, no. 1, February 1988.
- [7] J. S. Chase, H. M. Levy, M. Baker-Harvey, and E. D. Lazowska. "How to use a 64-bit virtual address space." Tech. Rep. 92-03-02, University of Washington, March 1992.
- [8] J. S. Chase, H. M. Levy, E. D. Lazowska, and M. Baker-Harvey. "Lightweight shared objects in a 64-bit operating system." Tech. Rep. 92-03-09, University of Washington, March 1992.
- [9] J. B. Chen, A. Borg, and N. P. Jouppi. "A simulation based study of TLB performance." In *Proc. 19th Annual International Symposium on Computer Architecture (ISCA-19)*, May 1992.
- [10] R. Cheng. "Virtual address cache in UNIX." In *Proceedings of the Summer 1987 USENIX Technical Conference*, June 1987.
- [11] D. R. Cheriton, H. A. Goosen, and P. D. Boyle. "Multi-level shared caching techniques for scalability in VMP-MC." In *Proc. 16th Annual International Symposium on Computer Architecture (ISCA-16)*, June 1989.
- [12] D. R. Cheriton, A. Gupta, P. D. Boyle, and H. A. Goosen. "The VMP multiprocessor: Initial experience, refinements and performance evaluation." In *Proc. 15th Annual International Symposium on Computer Architecture (ISCA-15)*, May 1988.
- [13] D. R. Cheriton, G. A. Slavenburg, and P. D. Boyle. "Software-controlled caches in the VMP multiprocessor." In *Proc. 13th Annual International Symposium on Computer Architecture (ISCA-13)*, January 1986.
- [14] D. R. Cheriton, G. A. Slavenburg, and P. D. Boyle. "Software-controlled caches in the VMP multiprocessor." In *Proceedings of the 1986 International Symposium on Computer Architecture*, June 1986.
- [15] H. Custer. "Inside Windows/NT." Tech. Rep., Microsoft Press, 1993.
- [16] H. Deitel. *Inside OS/2*. Addison-Wesley, Reading MA, 1990.
- [17] Digital. *DECchip 21064 and DECchip 21064A Alpha AXP Microprocessors Hardware Reference Manual*. Digital Equipment Corporation, Maynard MA, 1994.
- [18] P. Druschel and L. L. Peterson. "Fbufs: A high-bandwidth cross-domain transfer facility." In *Proc. Fourteenth ACM Symposium on Operating Systems Principles (SOSP-14)*, December 1993, pp. 189–202.
- [19] D. Engler, R. Dean, A. Forin, and R. Rashid. "The operating system as a secure programmable machine." In *Proc. 1994 European SIGOPS Workshop*, September 1994.
- [20] W. E. Garrett, M. L. Scott, R. Bianchini, L. I. Kontothanassis, R. A. McCallum, J. A. Thomas, R. Wisniewski, and S. Luk. "Linking shared segments." In *USENIX Technical Conference Proceedings*, January 1993.
- [21] W. E. Garrett, R. Bianchini, L. Kontothanassis, R. A. McCallum, J. Thomas, R. Wisniewski, and M. L. Scott. "Dynamic sharing and backward compatibility on 64-bit machines." Tech. Rep. TR 418, University of Rochester, April 1992.
- [22] J. R. Goodman. "Coherency for multiprocessor virtual address caches." In *Proc. Second Int'l Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS-2)*, October 1987, pp. 72–81.
- [23] J. Heinrich, Ed. *MIPS R10000 Microprocessor User's Manual, version 1.0*. MIPS Technologies, Inc., Mountain View CA, June 1995.
- [24] J. L. Hennessy and D. A. Patterson. *Computer Architecture: A Quantitative Approach*. Morgan Kaufmann Publishers, Inc., 1990.
- [25] Hewlett-Packard. *PA-RISC 1.1 Architecture and Instruction Set Reference Manual*. Hewlett-Packard Company, 1990.
- [26] M. D. Hill, S. J. Eggers, J. R. Larus, G. S. Taylor, G. Adams, B. K. Bose, G. A. Gibson, P. M. Hansen, J. Keller, S. I. Kong, C. G. Lee, D. Lee, J. M. Pendleton, S. A. Ritchie, D. A. Wood, B. G. Zorn,

- P. N. Hilfinger, D. Hodges, R. H. Katz, J. K. Ousterhout, and D. A. Patterson. "Design Decisions in SPUR." *IEEE Computer*, vol. 19, no. 11, November 1986.
- [27] J. Huck. *Personal communication*. 1996.
- [28] J. Huck and J. Hays. "Architectural support for translation table management in large address space machines." In *Proc. 20th Annual International Symposium on Computer Architecture (ISCA-20)*, May 1993.
- [29] IBM and Motorola. *PowerPC 601 RISC Microprocessor User's Manual*. IBM Microelectronics and Motorola, 1993.
- [30] J. Inouye, R. Konuru, J. Walpole, and B. Sears. "The effects of virtually addressed caches on virtual memory design and performance." Tech. Rep. CS/E 92-010, Oregon Graduate Institute, 1992.
- [31] Intel. *Pentium Processor User's Manual*. Intel Corporation, Mt. Prospect IL, 1993.
- [32] B. L. Jacob and T. N. Mudge. "Software-managed address translation." In *Proc. Third International Symposium on High Performance Computer Architecture (HPCA-3)*, San Antonio TX, February 1997.
- [33] G. Kane and J. Heinrich. *MIPS RISC Architecture*. Prentice-Hall, Englewood Cliffs NJ, 1992.
- [34] Y. A. Khalidi and M. Talluri. "Improving the address translation performance of widely shared pages." Tech. Rep. SMLI TR-95-38, Sun Microsystems, February 1995.
- [35] Y. A. Khalidi, M. Talluri, M. N. Nelson, and D. Williams. "Virtual memory support for multiple page sizes." In *Proc. Fourth Workshop on Workstation Operating Systems*, October 1993.
- [36] S. J. Leffler, M. K. McKusick, M. J. Karels, and J. S. Quarterman. *The Design and Implementation of the 4.3BSD UNIX Operating System*. Addison-Wesley Publishing Company, 1989.
- [37] J. Liedtke. "Improving IPC by kernel design." In *Proc. Fourteenth ACM Symposium on Operating Systems Principles (SOSP-14)*, December 1993, pp. 175–187.
- [38] J. Liedtke. "Address space sparsity and fine granularity." *ACM Operating Systems Review*, vol. 29, no. 1, pp. 87–90, January 1995.
- [39] J. Liedtke and K. Elphinstone. "Guarded page tables on MIPS R4600." *ACM Operating Systems Review*, vol. 30, no. 1, pp. 4–15, January 1996.
- [40] C. May, E. Silha, R. Simpson, and H. Warren, Eds. *The PowerPC Architecture: A Specification for a New Family of RISC Processors*. Morgan Kaufmann Publishers, San Francisco CA, 1994.
- [41] D. Nagle. *Personal communication*. 1995.
- [42] D. Nagle, R. Uhlig, T. Stanley, S. Sechrest, T. Mudge, and R. Brown. "Design tradeoffs for software-managed TLBs." In *Proc. 20th Annual International Symposium on Computer Architecture (ISCA-20)*, May 1993.
- [43] R. Rashid, A. Tevanian, M. Young, D. Young, R. Baron, D. Black, W. Bolosky, and J. Chew. "Machine-independent virtual memory management for paged uniprocessor and multiprocessor architectures." *IEEE Transactions on Computers*, vol. 37, no. 8, pp. 896–908, August 1988.
- [44] S. A. Ritchie. "TLB for free: In-cache address translation for a multiprocessor workstation." Tech. Rep. UCB/CSD 85/233, University of California, May 1985.
- [45] M. Rosenblum, E. Bugnion, S. A. Herrod, E. Witchel, and A. Gupta. "The impact of architectural trends on operating system performance." In *Proc. Fifteenth ACM Symposium on Operating Systems Principles (SOSP-15)*, December 1995.
- [46] M. L. Scott, T. J. LeBlanc, and B. D. Marsh. "Design rationale for Psyche, a general-purpose multiprocessor operating system." In *Proc. 1988 International Conference on Parallel Processing*, August 1988.
- [47] R. L. Sites, Ed. *Alpha Architecture Reference Manual*. Digital Equipment Corporation, Maynard MA, 1992.
- [48] M. Talluri and M. D. Hill. "Surpassing the TLB performance of superpages with less operating system support." In *Proc. Sixth Int'l Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS-6)*, San Jose CA, October 1994.
- [49] M. Talluri, M. D. Hill, and Y. A. Khalidi. "A new page table for 64-bit address spaces." In *Proc. Fifteenth ACM Symposium on Operating Systems Principles (SOSP-15)*, December 1995.
- [50] M. Talluri, S. Kong, M. D. Hill, and D. A. Patterson. "Tradeoffs in supporting two page sizes." In *Proc. 19th Annual International Symposium on Computer Architecture (ISCA-19)*, May 1992.
- [51] R. Wahbe, S. Lucco, T. E. Anderson, and S. L. Graham. "Efficient software-based fault isolation." In *Proc. Fourteenth ACM Symposium on Operating Systems Principles (SOSP-14)*, December 1993, pp. 203–216.
- [52] W.-H. Wang, J.-L. Baer, and H. M. Levy. "Organization and performance of a two-level virtual-real cache hierarchy." In *Proc. 16th Annual International Symposium on Computer Architecture (ISCA-16)*, June 1989, pp. 140–148.
- [53] D. L. Weaver and T. Germand, Eds. *The SPARC Architecture Manual version 9*. PTR Prentice Hall, Englewood Cliffs NJ, 1994.
- [54] S. Weiss and J. E. Smith. *POWER and PowerPC*. Morgan Kaufmann Publishers, San Francisco CA, 1994.
- [55] B. Wheeler and B. N. Bershad. "Consistency management for virtually indexed caches." In *Proc. Fifth Int'l Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS-5)*, October 1992.
- [56] D. A. Wood. *The Design and Evaluation of In-Cache Address Translation*. PhD thesis, University of California at Berkeley, March 1990.
- [57] D. A. Wood, S. J. Eggers, G. Gibson, M. D. Hill, J. M. Pendleton, S. A. Ritchie, G. S. Taylor, R. H. Katz, and D. A. Patterson. "An in-cache address translation mechanism." In *Proc. 13th Annual International Symposium on Computer Architecture (ISCA-13)*, January 1986.