

QuickStore: A High Performance Mapped Object Store

Seth J. White David J. DeWitt

Computer Sciences Department
University of Wisconsin
Madison, WI 53706
{white,dewitt}@cs.wisc.edu

ABSTRACT

This paper presents, QuickStore, a memory-mapped storage system for persistent C++ built on top of the EXODUS Storage Manager. QuickStore provides fast access to in-memory objects by allowing application programs to access objects via normal virtual memory pointers. The paper also presents the results of a detailed performance study using the OO7 benchmark. The study compares the performance of QuickStore with the latest implementation of the E programming language. These systems exemplify the two basic approaches (hardware and software) that have been used to implement persistence in object-oriented database systems. Both systems use the same underlying storage manager and compiler allowing us to make a truly apples-to-apples comparison of the hardware and software techniques.

1. Introduction

This paper presents, QuickStore, a memory-mapped storage system for persistent C++ built on top of the EXODUS Storage Manager (ESM) [Carey89a, Carey89b]. QuickStore uses standard virtual memory hardware to trigger the transfer of persistent data from secondary storage into main memory [Wilso90]. The advantage of this approach is that access to in-memory persistent objects is just as efficient as access to transient objects, i.e. application programs access objects by dereferencing normal virtual memory pointers, with no overhead for software residency checks as in [Moss90, Schuh90, White92].

QuickStore is implemented as a C++ class library that can be linked with an application, requiring no special compiler support. The memory-mapped architecture of QuickStore supports "persistence orthogonal to type", so that both transient and persistent objects can be manipulated using the same compiled code. Because QuickStore uses ESM to store persistent data on disk, it features a client-server architecture with full support for transactions (concurrency control and recovery), indices, and large objects. QuickStore places no additional limits on the size of a database, and the amount of data that can be accessed in the context of any single transaction is limited only by the size of virtual memory.

The paper also presents the results of a detailed performance study, in which we use the OO7 benchmark [Carey93] to compare the performance of QuickStore with the latest implementation of E [Rich89], a persistent programming language developed at Wisconsin that is also based on C++. The comparison between QuickStore and E is interesting because each of the systems takes a radically different approach toward implementing persistence.

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 of Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

SIGMOD 94- 5/94 Minneapolis, Minnesota, USA
© 1994 ACM 0-89791-639-5/94/0005..\$3.50

QuickStore employs a hardware faulting scheme that relies on virtual memory support (as mentioned above), while E uses an interpretive approach that is implemented in software.

These systems exemplify the two basic approaches (hardware and software) that have been used to implement persistence in object-oriented systems. Moreover, both QuickStore and E use the same underlying storage manager (ESM) and compiler. This allows us to make a truly apples-to-apples comparison of the hardware and software swizzling schemes, something which has not been done previously.

The remainder of the paper is organized as follows. Section 2 discusses related work on hardware and software based pointer swizzling schemes and points out how the performance results presented in this paper differ from previous studies. Section 3 describes the design of QuickStore. Section 4 presents our experimental methodology and Section 5 presents the results of the performance study. Section 6 contains some conclusions and proposals for future work.

2. Related Work

A detailed proposal advocating the use of virtual memory techniques to trigger the transfer of persistent objects from disk to main memory, first appeared in [Wilso90]. The basic approach described in [Wilso90] is termed "pointer swizzling at page fault time" since under this scheme all pointers on a page are converted from their disk format to normal virtual memory pointers (i.e. swizzled) by a page-fault handling routine before an application is given access to a newly resident page. In addition, pages of virtual memory are allocated for non-resident pages one step ahead of their actual use and access protected, so that references to these pages will cause a page-fault to be signaled. The technique described in [Wilso90] allows programs to access persistent objects by dereferencing standard virtual memory pointers, eliminating the need for software residency checks.

The basic ideas presented in [Wilso90] were, at the same time, independently used by the designers of ObjectStore [Objec90, Lamb91], a commercial OODBMS product from Object Design, Inc. The implementation of ObjectStore, outlined briefly in [Objec90], differs in some interesting ways from the scheme described in [Wilso90]; most notably in the way that pointer swizzling is implemented, and in how pointers are represented on disk.

Under the approach outlined in [Objec90], pointers between persistent objects are stored on disk as virtual memory pointers instead of being stored in a different disk format as in [Wilso90]. In other words, pointer fields in objects simply contain the value that they last were assigned when the page was resident in main memory in ObjectStore. When a page containing persistent objects is first referenced by an application program, ObjectStore attempts to assign the page to the same virtual address as when the page was

This research is sponsored by the Advanced Research Project Agency, ARPA order number 018 (formerly 8230), and monitored by the U.S. Army Research Laboratory under contract DAAB07-91-C-Q518.

last memory resident. If all of the pages accessed by an application can be assigned to their previous locations in memory, then the pointers contained on the pages can retain their previous values, and need not be "swizzled", i.e. changed to reflect some new assignment of pages to memory locations, as part of the faulting process. If any page cannot be assigned to its previous address (because of a conflict with another page), then pointers that reference objects on the page will need to be altered (i.e swizzled) to reflect the new location of the page.

This scheme requires that the system maintain some additional information describing the previous assignment of disk pages to virtual memory addresses. The hope is that processing this information will be less expensive on average, than swizzling the pointers on pages that are faulted into memory by the application program. We note here that QuickStore is similar to ObjectStore in that QuickStore also stores pointers on disk as virtual memory pointers. Section 3 contains a detailed discussion of the implementation of QuickStore.

The Texas [Singh92] and Cricket [Shek90] storage systems also use virtual memory techniques to implement persistence. Texas stores pointers on disk as 8-byte file offsets, and swizzles pointers to virtual addresses as described in [Wilso90] at fault time. Currently, all data is stored in a single file (implemented on a raw *Unix* disk partition) in Texas [Singh92]. Although, QuickStore and Texas are different in their implementation details, there are some similarities between the two systems. For example, both systems are implemented as C++ libraries that add persistence to C++ programs without the need for compiler support. Both systems also support the notion of "persistence orthogonal to type". This allows the same compiled code to manipulate both transient and persistent objects. Both systems also allow the database size to be bigger than the size of virtual memory.

Texas, however, is currently a single user, single processor system while QuickStore, since it is built on top of client-server EXODUS, features a client-server architecture with full transaction support including concurrency control, recovery, and support for distributed transactions. QuickStore is also different in that it manipulates objects directly in the ESM client buffer pool, while Texas copies objects into a separate heap area allocated in virtual memory. This limits the amount of data that can currently be accessed during a single transaction by Texas to the size of the disk swap area backing the application process. QuickStore also manages paging in the ESM client buffer pool explicitly, while Texas simply allows pages to be swapped to disk by the virtual memory subsystem when the process size exceeds the size of physical memory. Cricket, on the other hand, uses the Mach external pager facility to map persistent data into an application's address space (see [Shek90] for details).

We next discuss previous performance studies of pointer swizzling and object faulting techniques, and point out how the study presented here differs from them. [Moss90] contains a study of several software swizzling techniques and examines various issues relevant to pointer swizzling. Among these are whether swizzling has better performance than simply using object identifiers to locate objects, and whether objects should be manipulated in the buffer pool of the underlying storage manager, or copied out into a separate area of memory before swizzling takes place. [Moss90] also looks at lazy vs. eager swizzling. Eager swizzling involves prefetching the entire collection of objects into memory so that all pointers can be swizzled, while lazy swizzling swizzles pointers incrementally as objects are accessed and faulted into memory by the application program.

We do not consider copy swizzling approaches since [White92] showed that they do not perform well when the database size is larger than physical memory. The study presented here also differs from [Moss90] in that we allow pages of objects to be replaced in

the buffer pool, while [Moss90] only considers small data sets where no paging occurs. The systems we examine also include concurrency control and recovery, while those examined in [Moss90] did not.

[Hoski93] examines the performance of several object faulting schemes in the context of a persistent Smalltalk implementation. [Hoski93] includes one scheme that uses virtual memory techniques to detect accesses to non-resident objects. The approach described in [Hoski93] allocates fault-blocks, special objects that stand in for non-resident objects, in protected pages. When the application tries to access an object through its corresponding fault block, an access violation is signaled. The results presented in [Hoski93] show this scheme to have very poor performance. It is not clear, however, whether this is due to the overhead associated with using virtual memory or is the result of extra work that must be performed during each object fault to locate and eliminate any outstanding pointers to the fault block that caused the fault. This work involves examining the pointer fields of all transient and persistent objects that contain pointers to the fault block. Finally, we note that the effects of page replacement in the buffer pool and updates are also not considered in [Hoski93].

In [White92] the performance of several implementations of the E language [Rich89, Rich90, Schuh90] and ObjectStore [Objec90, Lamb91], a commercial OODBMS, are compared. The results presented in [White92] were inconclusive, however, in providing a true comparison of software and hardware-based schemes since the underlying storage managers used by the systems were different and because the systems used different compilers. In the study presented in this paper, all of the systems use the same underlying storage manager and compiler, so any differences in performance are due to the swizzling and faulting technique that was used.

One additional difference between the systems compared in [White92] and those examined here, is that the systems included in the current study are much less restrictive in terms of the amount of data that can be accessed during a transaction, and all systems manage paging of persistent data explicitly. This differs from the approach used by EPVM 2.0 in [White92], which limited the amount of data that could be accessed during a transaction to the size of the disk swap area backing the process, and which allowed objects to be swapped to disk by the virtual memory subsystem when the size of the process exceeded the size of physical memory.

3. QuickStore Design Concepts

3.1. Overview of the Memory-Mapped Architecture

As mentioned in Section 1, QuickStore uses ESM to store persistent objects on disk. ESM features a page-shipping architecture, in which objects are transferred from the server to the client a page-at-a-time. Once a page of objects has been read into the buffer pool of the ESM client, applications that use QuickStore access objects on the page directly in the ESM client buffer pool, by dereferencing normal virtual memory pointers. Objects are always accessed in the context of a transaction in QuickStore.

To understand the way that QuickStore coordinates access to persistent objects, it is useful to view the virtual address space of the application process as being divided into a contiguous sequence of *frames* of equal length. In our case, these frames are 8 K-bytes in size, the same size as pages on disk. The ESM client buffer pool can also be viewed as a (much smaller) sequence of 8 K-byte frames. To coordinate access to persistent objects, QuickStore maintains a physical mapping from virtual memory frames to frames in the buffer pool. This physical mapping is **dynamic**, since paging in the buffer pool requires that the same frame of virtual memory be mapped to different frames in the buffer pool at different points in time. The mapping can also be viewed as a

logical mapping from virtual memory frames to disk pages. When viewed this way, the mapping is *static* since the same virtual frame is always associated with the same disk page during a transaction.

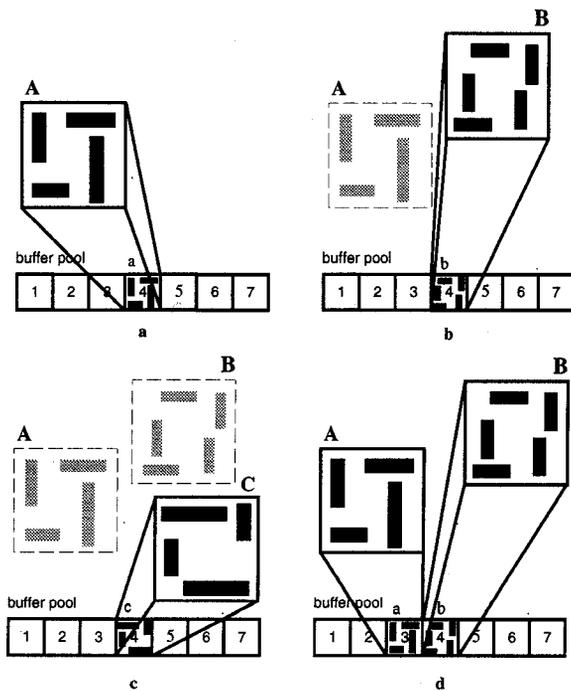


Figure 1. Mapping virtual frames into the buffer pool.

Figure 1 illustrates this mapping scheme in more detail. The buffer pool shown in Figure 1 contains 7 frames (labeled from 1 to 7). Virtual memory frames are denoted using upper-case letters, while disk pages are specified in lower-case. In the discussion that follows, we sometimes refer to the virtual memory frame beginning at address *A* as frame *A*.

The virtual memory frame corresponding to a disk page that contains persistent data is selected by QuickStore and access protected, before the page can be accessed by an application. When the application first attempts to access an object on the page by dereferencing a pointer into its frame, a page-fault is signaled and a fault handling routine that is part of the QuickStore runtime system is invoked. This fault handling routine is responsible for reading the page from disk, updating various data structures, and enabling access permission on the virtual frame that caused the fault so that execution of the program can resume. For example, in Figure 1a page *a* has been read from disk into frame 4 of the buffer pool. Page *a* is "mapped" to virtual address *A*. Read access has been enabled on frame *A*, so that the application can read the objects contained on page *a*. We note that once the mapping from virtual address *A* to page *a* has been established, the application can access objects on page *a* by dereferencing pointers to frame *A* at any time. Thus, the mapping from *A* to *a* must remain valid until the end of the current transaction (or longer, if requested) in order to preserve the semantics of any pointers that the application may have to objects on page *a*.

If the objects on page *a* contain pointers to objects on other non-resident pages, then virtual frames are assigned to these pages when page *a* is faulted into memory, if they haven't been already. A frame for a non-resident page remains access protected until the program attempts to dereference a pointer into the frame. Figure 1 doesn't explicitly show any frames of this type for page *a* since they are not important to the current discussion.

When the buffer pool becomes full, paging will occur and page *a* may be selected for replacement by the buffer manager. This is what has happened in Figure 1b. Here, page *b* has been read from disk into frame 4 of the buffer pool, replacing page *a*. Page *b* has been mapped to virtual address *B* and read access on *B* has been enabled. Note that since we assume that the buffer pool is full in Figure 1b, additional virtual frames (not shown) will also have been mapped to the remaining 6 frames in the buffer pool other than frame 4. If the application continues to access additional pages of objects in the database, then the situation shown in Figure 1c may result. In Figure 1c, page *c* has been read into memory and replaced page *b* in frame 4 of the buffer pool. Page *c* has been mapped to virtual frame *C* and read access on *C* has been enabled. This illustrates that, in general, any number of virtual frames may be associated with a particular frame of the buffer pool over the course of a transaction.

The reader may be wondering at this point, what would happen in Figure 1b if the application attempted to dereference pointers into virtual frame *A* after page *a* has been replaced in the buffer pool by page *b*? Won't these pointers refer to data on page *b*? This problem is avoided by disabling read access on frame *A* when page *a* is not in memory. If the application again dereferences pointers into frame *A*, a page-fault will be signaled and the fault handling routine invoked. The fault handling routine will call ESM to reread page *a*, map virtual frame *A* to the frame in the buffer pool that now contains *a*, and enable read permission on frame *A* once again.

To illustrate this, Figure 1d shows what might result if page *a* were immediately referenced after it was replaced in Figure 1b. In this case, *a* has been reread by ESM into frame 3 in the buffer pool and frame 3 has been mapped to virtual memory address *A*. This further illustrates the dynamic nature of the physical mapping from virtual memory frames to frames in the buffer pool since virtual frame *A* is mapped to buffer frame 4 in Figure 1a and remapped to buffer frame 3 in Figure 1d. However, the mapping between virtual frames and disk pages is static since virtual frame *A* is always mapped to disk page *a*.

3.2. Implementation Details

QuickStore uses the UNIX *mmap* system call to implement the physical mapping from virtual memory frames to frames in the ESM client buffer pool, and to control virtual frames' access protections. It was necessary to modify the ESM client software slightly in order to accommodate the use of *mmap* since *mmap* really just associates virtual memory addresses with offsets in a file, while ESM normally calls the UNIX function *malloc* to allocate space in memory for its client buffer pool. To make ESM and *mmap* work together, the buffer pool allocation code was changed so that it would first open a file (and resize it if necessary) equal in size to the size of the client buffer pool. Then the buffer allocation code calls *mmap* to associate a range of virtual memory with the entire file. The rest of the ESM client software uses this range of memory to access the buffer pool just as though the memory had been allocated using *malloc*.

The important thing to note is that the file serves as backing store for the buffer pool. Swap space and actual physical memory are never allocated for the virtual frames that are mapped into the file by *mmap*, so mapping a huge amount of virtual memory into the buffer pool doesn't affect the size of the process, although it may increase the size of page tables maintained by the operating system. One should also note that the contiguous range of addresses used by the ESM client to access the buffer pool is different from the 8 K-byte ranges of addresses that the application program uses to access pages in the buffer pool. The former is simply used to integrate an already existing storage manager (ESM) with the memory mapped approach and would not, in general, be required by a memory mapped implementation.

We would like to point out, however, that using *mmap* in the way that we did, actually caused some minor performance problems in the implementation. Because the workstation used as the client machine (a Sun ELC) in the benchmark experiments had a virtually mapped CPU cache, accessing the same page of physical memory in the buffer pool via different virtual address ranges caused the CPU cache to be flushed whenever the process switched between the address ranges. This increased the number of *min faults*, virtual memory page faults that do not require I/O in Unix terminology, experienced by the application. We note the effects of this phenomena when discussing the performance results (see Section 5).

3.3. In-Memory Data Structures

QuickStore maintains an in-memory table that keeps track of the current logical mapping from virtual memory frames to disk pages. At a given point in time, the table contains an entry for every page that has been faulted into memory, plus entries for any additional pages which are referenced by pointers on these pages. Entries in the table are called *page descriptors* and are 60 bytes long. We note that disk pages themselves come in two types: pages that contain sets of objects that are smaller than a disk page are called *small object pages*, while pages that contain individual pages of multi-page objects are called *large object pages*. Table entries for small object pages and large object pages differ in some respects, so they are discussed separately.

A page descriptor for a small object page contains the range of virtual addresses associated with the page, the physical address of the page on disk, and a pointer to the page when it is pinned in the buffer pool. The physical address of the page, in our implementation, is the OID of a special meta-object (24 bytes) located on each small object page. Page descriptors also contain other fields such as flags that indicate what types of access are currently allowed on the frame associated with the page (read, write, and none), whether an exclusive lock has been obtained, and whether or not the page has previously been read into memory during the current transaction. This last flag is useful since it is not necessary to do any swizzling work for a page when it is reread during a transaction, since the pointers on the page are guaranteed to be valid.

The scheme used for large object pages is somewhat more complicated than the scheme for small object pages. The virtual memory frames associated with a multi-page object must be contiguous, so they are reserved all at once. To avoid maintaining individual table entries for every page of a multi-page object, multi-page objects that have not been accessed, but which are in the mapping, are represented by a single entry in the mapping table. The range of virtual addresses in this entry is the entire range of contiguous addresses associated with the object and the physical address field contains the OID of the object. When the first page of a multi-page object is accessed by the application program, the table entry is split so that there is one entry in the table for the page that has been accessed, and an entry for each contiguous sub-sequence of unaccessed pages. Table entries for sub-sequences of unaccessed pages of a multi-page object are split in turn when one of the pages contained in the sub-sequence is accessed.

The table organizes page descriptors according to the range of virtual memory addresses that they contain using a height balanced binary tree. One reason for using a binary tree was that it makes the splitting operation associated with large objects efficient. It is also helpful to keep the ranges of addresses currently allocated to persistent data ordered. For example, our current scheme for allocating virtual frames to disk pages uses a global counter (stored on disk) that is incremented by the frame size each time that a frame is allocated to a disk page. If the database becomes bigger than the size of virtual memory then this counter will wrap around and it may become necessary to scan the in-memory binary tree in order

to find a virtual frame that is currently not in use.

Page descriptors are also hashed based on their physical address (OID) and inserted into a hash table. (For large objects only the page descriptor that the beginning subsequence of the object is inserted into the hash table.) The hash table implements a reverse mapping from physical disk address to virtual memory address. The hash table is used by the fault handling routine as part of the pointer swizzling process (see below for details).

3.4. Pointer Swizzling in QuickStore

Like ObjectStore [Objec90, Lamb91], QuickStore stores pointers on disk as virtual memory addresses in exactly the same format that they have when they are in memory. Since virtual memory pointers are only meaningful in the context of an individual process, this scheme requires that the system maintain some additional meta-data that associates pointers on disk with the objects that they reference. The remainder of this section describes how this meta-data is stored in QuickStore.

QuickStore associates meta-data with individual disk pages. In the case of small object pages, each page contains a direct pointer (OID) to a *mapping object* containing the meta-data for the page. (Actually, the pointer is contained in a special meta-object located on the page.) The term *mapping object* is used since the object records the mapping between virtual frames referenced by pointers on the page and disk pages that was in effect when the page was last memory resident. Mapping objects are essentially just arrays of <virtual address frame, disk address> pairs.

Mapping information is stored separately instead of on the disk pages containing objects themselves because the space required to store the mapping information for a page can vary over time. For example, if the pointers on a page are updated, the number of frames referenced by pointers on the page may change, changing the number of entries in the mapping object. Multi-page objects are implemented similarly to small object pages, except that there is an array of meta-objects appended to the end of the large object containing one meta-object for each page of the large object. Finally, we note that each meta-object also contains a pointer (OID) to a bitmap object that records the locations of pointers on the page so that they can be swizzled. QuickStore uses a modified version of *gdb* to get the type information for objects that is used to maintain the bitmaps associated with pages.

To illustrate how the structures mentioned above are used, consider the actions that are taken when a page containing data is first read into memory by QuickStore. After reading the page, the meta-object on the page is examined and the OID of the mapping object that it contains is used to read the mapping object itself from disk. The runtime system then looks up the disk address contained in each mapping object entry in the in-memory table to see if the disk page is currently part of the mapping. If no entry is found in the table, then one is created using the information contained in the mapping object. The disk page will be assigned to its previous virtual frame at this point, if it is unused, or else a new frame is selected. If an entry for the disk page is found in the table, the system checks to see if the page is currently associated with the same virtual frame as the one in the mapping object entry.

If all of the disk pages in the mapping object are associated with their old virtual frames, then the swizzling process terminates. If some disk pages have been mapped to new locations, however, then the bitmap object is read from disk and used to find and update any pointers on the page that reference these pages. Note that even though bitmap objects are fixed in size, they are stored separately from their corresponding data page since they hopefully will not have to be read in most cases.

3.5. Buffer Pool Management

Most buffer managers in traditional database systems have used a clock style algorithm to approximate an LRU page replacement policy. We also felt that a clock algorithm was the best choice for use in QuickStore, however, implementing this type of scheme turned out to be more difficult in the context of a memory-mapped system where objects in the buffer pool are accessed by dereferencing virtual memory pointers. The reason for this is that there is less information available to the buffer manager indicating which pages have been accessed recently.

Recall that in traditional implementations of clock, a bit is usually kept for each frame in the buffer pool, indicating whether or not the frame has been accessed since the clock hand last swept over it. This bit is set by the database system each time the page is accessed and reset by the clock algorithm. There is no way to set such a flag, however, when dereferencing a pointer as in QuickStore.

One solution to this problem is to have the clock algorithm access-protect the virtual frame corresponding to a buffer pool frame when the clock hand reaches it. If the frame is subsequently reaccessed a page-fault will occur and the fault handling routine can re-enable access to the page. This scheme replaces the usual setting and unsetting of bits in a traditional clock algorithm with the enabling and disabling of access permissions on virtual memory frames. We experimented with this solution, however, in our experience the extra overhead of manipulating the page protections and handling additional page-faults made this approach prohibitively expensive in terms of performance.

To avoid the problem described above, QuickStore uses a simplified clock algorithm. Under this scheme the clock hand begins its sweep from wherever it stopped during the previous invocation of the clock algorithm. As soon as the clock hand reaches a page for which access is not enabled the algorithm selects that page for replacement. If the clock hand reaches the end of the buffer pool without finding a candidate for replacement, however, then the **entire** virtual address space of the process being used for persistent data is reprotected with a single call to *mmap* and the algorithm is restarted. This scheme performed much better than the original scheme outlined above in our experiments, and compared favorably with the more traditional clock replacement algorithm used by E (see Section 5).

3.6. Recovery

Implementing recovery for updates poses some special problems in the context of a memory-mapped scheme as well. For example, since application programs are able to update objects by dereferencing virtual memory pointers, it is difficult to know what portions of objects have been modified and require logging. Furthermore, it is desirable to batch the effects of updates together and log them all at once, if possible, since some applications may update the same object many times during a transaction.

Due to the considerations mentioned above, we decided to use a page diffing scheme to generate log records for objects that have been updated in QuickStore. Virtual memory frames that are mapped to pages in the database that have not been updated, do not have write access enabled, so the first attempt by the application program to update an object on a page will cause a page-fault. When the fault-handler detects that an access violation is due to a write attempt, it copies the original values contained in the objects on the page into an in-memory heap data structure. The fault-handler also obtains an exclusive lock on the page from ESM, if needed, and enables write access on the virtual frame that caused the fault before returning control to the application. The application program can then update the objects on the page directly in the buffer pool.

At transaction commit time, or sooner if paging in the buffer pool occurs or the heap becomes full, the old values of objects contained in the heap and the corresponding updated values of objects in the buffer pool are diffed to determine if log records need to be generated. The processes of diffing and generating log records are interleaved in QuickStore. To understand why this is the so, consider as an example the case when the first and last byte of an 1 K-byte object have been updated. In this case QS minimizes the amount of data written to the log by generating two log records, one for each modified byte, instead of one big log record for the entire object. On the other hand if the first, third, and fifth bytes had been modified, QS would generate a single log record for the first five bytes of the object. This is cheaper than generating multiple log records since each log record contains a relatively large (~50 byte) header area for storing information needed by the ESM recovery scheme.

Care must also be taken when processing updates to update the mapping tables associated with each modified page if necessary. Recall that the mapping table for a page keeps track of the set of pages that are referred to by pointers on the page. Updates to objects on a page can change the pages that are members of this set, making it necessary to update the mapping tables as well. Updating the mapping table for a page requires that each pointer contained on the page be examined and the in-memory table consulted to determine which page in the database it references. The bitmap for the page is used to locate the pointers that it contains and from these pointers a new set of referenced pages is constructed. This new set is then compared element by element with the old set to see if the set has changed. If it has, then the mapping object for the page is updated to reflect the new set of referenced pages.

4. Performance Experiments

This section briefly describes the structure of the OO7 benchmark database and the benchmark operations that were included in the performance study. (See [Carey93] for a complete description of the OO7 benchmark.) The hardware and software systems included in the study are also discussed.

4.1. The OO7 Benchmark Database

The OO7 database is intended to be suggestive of many different CAD/CAM/CASE applications. We used two sizes for the OO7 database: small and medium. Table 1 summarizes the parameters of the database. A key component of the database is a set of *composite parts*. Each composite part is intended to suggest a design primitive such as a register cell in a VLSI CAD application. The number of composite parts was set to 500 in both the small and

Parameter	Small	Medium
NumAtomicPerComp	20	200
NumConnPerAtomic	3	3
DocumentSize (bytes)	2000	20000
Manual Size (bytes)	100K	1M
NumCompPerModule	500	500
NumAssmPerAssm	3	3
NumAssmLevels	7	7
NumCompPerAssm	3	3
NumModules	1	1

Table 1. OO7 Benchmark database parameters.

medium databases. Each composite part has a number of attributes, including the integer attributes **id** and **buildDate**. Associated with each composite part is a *document* object that models a small amount of documentation associated with the composite part. Each document has an integer attribute **id**, a small character attribute **title** and a character string attribute **text**. The length of the string attribute is controlled by the parameter *DocumentSize*.

Each composite part also has an associated graph of *atomic parts*. Intuitively, the atomic parts within a composite part are the units out of which the composite part is constructed. One atomic part in each composite part's graph is designated as the "root part". In the small database, each composite part's graph contains 20 atomic parts, while in the medium database, each composite part's graph contains 200 atomic parts. Each atomic part has the integer attributes **id**, **buildDate**, **x**, **y**, and **docId**. Each atomic part is connected via a bi-directional association to *NumConnPerAtomic* other atomic parts which was set to three in the experiments. The connections between atomic parts are implemented by interposing a connection object between each pair of connected atomic parts.

Additional structure is imposed on the set of composite parts by the "assembly hierarchy". Each assembly is either made up of composite parts (in which case it is a *base assembly*) or it is made up of other assembly objects (in which case it is a *complex assembly*). The first level of the assembly hierarchy consists of *base assembly* objects. Base assembly objects have the integer attributes **id** and **buildDate**. Each base assembly has a bi-directional association with three composite parts which are chosen at random from the set of all composite parts. Higher levels in the assembly hierarchy are made up of *complex assemblies*. Each complex assembly has a bi-directional association with three subassemblies, which can either be base assemblies (if the complex assembly is at level two in the assembly hierarchy) or other complex assemblies (if the complex assembly is higher in the hierarchy). There are seven levels in the assembly hierarchy.

Each assembly hierarchy is called a *module*. Modules are intended to model the largest subunits of the database application. Modules have several scalar attributes. Each module also has an associated *Manual* object, which is a larger version of a document. Manuals are included for use in testing the handling of very large (but simple) objects. Figure 2 depicts the full structure of the single user OO7 Benchmark Database.

4.2. The OO7 Benchmark Operations

This section describes the OO7 benchmark operations that were used in the study. The full set of benchmark operations, consists of a set of 10 tests termed traversals, and another set of 8 query tests. We do not present results for the queries since none of the systems we tested has a declarative query language. Some of the traversal operations were also omitted because they didn't highlight any additional differences among the systems that were studied. The traversal tests are numbered from one to ten.

The T1 traversal performs a depth-first traversal of the assembly hierarchy. As each base assembly is visited, each of its composite parts is visited and a depth first search on the graph of atomic parts is performed. The traversal returns a count of the number of atomic parts visited, but otherwise no additional work is performed. The T6 traversal is similar to T1, except that instead of visiting the entire graph of atomic parts for each composite part, T6 just visits the root atomic part.

T2 and T3 are also similar to T1, but they add updates. Each T2 traversal increments the (x,y) attributes contained in atomic parts as follows¹:

- T2A—Update the root atomic part of each composite part.
- T2B—Update all atomic parts of each composite part.
- T2C—Update all atomic parts four times.

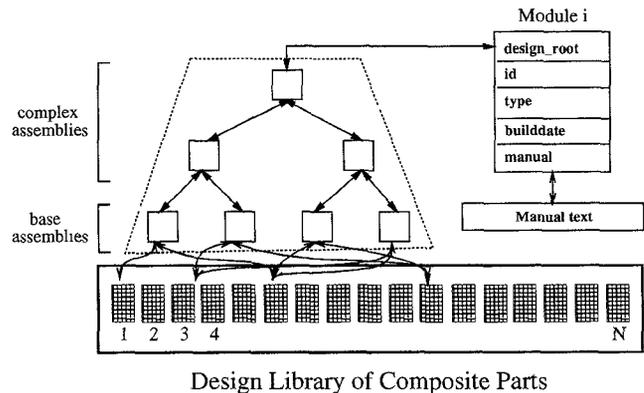


Figure 2. Structure of a module.

The T3 traversals are similar to T2 except that the **buildDate** field of atomic parts is incremented. This field is indexed, so T3 highlights the cost of updates of indexed fields. We used 3 traversals that are not based on T1. T7 picks a random atomic part and traverses up to the root of the design hierarchy. T8 scans the manual object associated with the module and counts the occurrences of a specified character, and T9 compares the first and last characters of the manual to see if they are equal.

4.3. Systems Tested

4.3.1. E

This section briefly describes the current implementation of the E language. E and QuickStore both offer basically the same functionality, however, E implements persistence using a software interpreter, EPVM 3.0. EPVM 3.0 has a functional interface, so operations such as dereferencing an unswizzled pointer in E are handled by calling an EPVM function to perform the dereference. As part of handling a reference to a persistent object, EPVM may in turn call ESM if the page containing an object is not in memory, and update its own internal data structures before returning control to the application. In addition to calls of EPVM functions, the code generated by the E compiler (a modified version of the gnu C++ compiler) contains in-line code sequences to handle certain basic operations. For example, residency checks and dereferences of swizzled pointers are done in-line and do not require a function call, which improves performance.

Like QuickStore, EPVM 3.0 accesses memory-resident persistent objects directly in the ESM client buffer pool. The interpreter maintains a hash table that contains an entry for each page of objects that is currently in memory. The pointer swizzling scheme used in EPVM 3.0 is similar to the scheme used in EPVM 1.0 [Schuh90] except that swizzled pointers point directly to objects in the buffer pool. This swizzling scheme only swizzles pointers that are local variables in C++ functions. Pointers within persistent objects are not swizzled because this makes page replacement in the buffer pool difficult [White92].

¹[Carey93] specifies that the (x, y) attributes should be swapped, however, we increment them instead, so that multiple updates of the same object always change the object's value. This guarantees that the diffing scheme used by QuickStore for recovery will always generate a log record.

Update operations on persistent objects are always handled by an interpreter function in EPVM 3.0. The update scheme used copies the original values of objects into a side buffer, and updates the objects in place in the buffer pool. Original values of objects and updated values are used to generate log records at transaction commit, or sooner if the side buffer or the buffer pool become full. However, no diffing is performed as in QuickStore. EPVM 3.0 employs a scheme that breaks large objects into 1K chunks for logging purposes. Objects that are smaller than 1K are logged in their entirety.

4.3.2. QuickStore

A detailed description of QuickStore is given in Section 3. Although QuickStore and E offer nearly the same functionality, it is important to point out one fundamental way in which the two systems differ. This has to do with the support that both systems provide for the notion of object identity. Section 3 described the scheme used by QuickStore (and ObjectStore [Objec90]) to implement a mapping between virtual memory frames and disk pages. This mapping is maintained for pages when they are in memory as well as when they reside on disk. Because of this mapping, pointers to persistent objects can be viewed as a <virtual frame, offset> pair where the high order bits of the pointer identify the virtual memory frame referenced by the pointer and the low order bits specify an offset into that frame. Virtual memory frames are mapped to disk pages, so pointers really just specify offsets or locations on pages.

To see why this scheme doesn't support object identity, consider what happens when an object, for which there are outstanding references, is deleted. The page that contained the object can be faulted into memory by subsequent program runs (assuming there are other objects on the page) and mapped to some virtual memory frame. If the program then dereferences dangling pointers to the deleted object, no error will be explicitly flagged. If a new object occupies the space on the page previously used by the deleted object then the dangling pointers will reference this object.

QuickStore doesn't fully support object identity or "checked references" to objects because the overhead would be prohibitive. For example, the meta-data that would be required to associate every unique pointer on a page with its corresponding OID would likely be an order of magnitude greater than the current scheme used by QuickStore. Furthermore, we are aware of no commercial or research system (including ObjectStore) that supports "checked references" for normal pointers in the context of a memory-mapped scheme. E, on the other hand, supports object identity fully, including checked references. E implements this by storing pointers as full 12 byte OIDs within objects. This is a reasonable approach, but it does incur certain costs. For example, since objects are larger in E than in QuickStore the database as a whole is larger, and E generally performs more I/O. Also, dereferencing big pointers is more expensive in terms of CPU requirements than dereferencing virtual memory pointers.

Because of these differences, we included a third system in the performance study. This system is identical to QuickStore, except that each object in the database has been padded so that it is the same size as the corresponding object in E. Comparing the performance of this system to the performance of E in the experiments where faults take place gives insight into the overhead of faulting for the memory-mapped approach, while comparing it with QuickStore indicates the advantage gained by QuickStore due to its smaller object size. In addition, one can think of this system as approximating the performance of a hybrid memory-mapped scheme that allows large pointers to be embedded within objects, thus supporting both checked and unchecked references.

4.4. Hardware and Software Used

As a test vehicle we used a pair of Sun workstations on an isolated Ethernet. A Sun IPX workstation configured with 48 megabytes of memory, one 424 megabyte disk drive (model Sun0424) and one 1.3 gigabyte disk drive (model Sun1.3G) was used as the server. The Sun 1.3G drive was used by ESM to hold the database, and the second Sun 0424 drive was used to hold the ESM transaction log. The data and recovery disks were configured as raw disks. For the client we used a Sun Sparc ELC workstation (about 20MIPS) configured with 24 megabytes of memory.

The systems included in the study used the client-server version of EXODUS (ESM V3.0). During the experiments ESM used a disk page size of 8 Kbytes (this is also the unit of transfer between a client and the server). The client and server buffer pools were set to 1,536 (12 MBytes) and 4,608 pages (36 Mbytes) respectively. Release 4.1.3 of the SunOS was run on both workstations used in the experiments. QuickStore was compiled using the GNU g++ compiler V2.3.1. The E compiler is a modified version of the GNU compiler.

5. Performance Results

5.1. Database Sizes

The size of the OO7 database is important in understanding the performance results. Table 2 shows the database sizes for E, QuickStore (QS), and QuickStore with big objects (QS-B)². The QS database is roughly 60% as big as the E database for both the small and medium cases. This is because of the different schemes used by the systems to store pointers. The QS-B database is bigger than the E database due to the overhead for storing bitmaps that indicate the locations of pointers on pages and mapping tables.

	Small	Medium
QS	6.6	54.2
E	10.5	94.1
QS-B	11.5	98.5

Table 2. Database Sizes (in megabytes)

5.2. Small Cold Results

This section presents the cold results for the small database experiments. The cold results were obtained by running the OO7 benchmark operations when no data was cached in memory at either the client or server machines. The times presented represent the average of 10 runs of the benchmark operations, except where noted otherwise. The times were computed by calling the Unix function *gettimeofday* which had a granularity of several microseconds on the client machine.

Figure 3 presents the cold response times for the read-only traversals included in the study. Table 3 gives the number of client I/O requests. As Figure 3 shows, QS is 37% faster than E during T1³. This difference in performance is largely due to the smaller database size for QS which causes it to read 53% fewer pages from disk than E (Table 3). The overwhelming majority of I/O activity during T1 is due to reading clusters of composite parts. Each composite part cluster occupied a little less than one page for QS, while

²Objects in QS-B are padded to the same size as the corresponding objects in the E implementation.

³T1: DFS of assembly hierarchy visiting all atomic parts.

for E close to two disk pages were required. This accounts for the roughly 2 to 1 ratio in the number of disk reads between the two systems. Comparing E with QS-B, shows that QS-B is 15% slower than E during T1. QS-B always issues slightly more I/O requests than E since QS-B must also read mapping tables to support the memory-mapping scheme.

The performance of QS is only 4% better than E during T6⁴. As in T1, differences in the size of composite part clusters between the two systems play an important role in determining their relative performance. Table 3 shows that the amount of I/O for QS is almost the same during both T1 and T6, while the number of disk reads for E decreases by 41% during T6. E does noticeably fewer I/O operations during T6 because it generally doesn't read the entire composite part cluster as QS does. The performance of QS-B is 27% slower than E during T6. As the detailed faulting times (shown below) will illustrate, this difference in performance is close to the actual percentage difference in individual page fault costs for the systems, as CPU costs have less of an overall impact on performance during T6 than during T1.

E has the best performance during T7⁵. QS is 20% slower than E because of increased faulting costs relative to T1 and T6. Faults are more expensive for QS during T7 because a large fraction of faults (86%) are spent reading pages of *assembly* objects. These pages have larger mapping tables because pointers from *association* objects to *composite part* objects are uniformly distributed among all *composite part* objects in the database. This increases the average I/O cost for reading the mapping tables and the number of table entries (139 on average for T7 vs 20 on average for T1) that must be examined per fault. QS-B is 34% slower overall

during T7 relative to E.

Turning to T8⁶, Figure 3 shows that E is roughly 3 times slower than QS. This is because the E interpreter is invoked once for each character of the manual that is examined by T8, while QS simply has to dereference a virtual memory pointer to access each character. By contrast, Figure 3 shows that E is nearly twice as fast as QS on T9⁷. This difference is due almost entirely to faulting costs since very little work is done on the objects faulted in during T9. It is not surprising that faulting costs for QS are relatively high in this case since T9 touches very few pages. QS and QS-B have similar performance during T8 and T9 since character data is the same size for both systems.

The results presented in Figure 3 for the read-only traversals have demonstrated that the per page faulting cost for the memory mapped approach is higher than for the software approach. For example, QS-B almost always has slower performance than E. To determine the magnitude of this difference, we next examine the individual faulting costs of the systems in more detail. Table 4 shows the average cost per fault in milliseconds for each of the systems during T1 and T6. These times were calculated by subtracting the time required to execute a hot traversal from the time required for a cold traversal, and then dividing the result by the number of page faults to get the average per fault cost. To make sure that the results obtained were accurate, the numbers used to perform the calculation represented the average of 100 runs of each traversal experiment.

According to Table 4, the faulting cost for QS-B is slightly higher than for QS. We speculated that this was because the mapping tables for QS-B were larger than in QS. However, this turned out not to be the case since the average number of mapping table entries in the small database was 16 for QS and only 12 for QS-B. The comparison between QS and E in Table 4 is more interesting. It shows that individual page faults are roughly 20% more expensive for QS during the T1 and T6 traversals. The corresponding figure for QS-B and E averages 26%, which correlates closely with the difference in response time between QS-B and E during T6.

To better understand the additional faulting overhead of the memory mapped scheme, Table 5 shows a detailed breakdown of the average faulting time for QS. As a check we present detailed numbers for both T1 and T6. One would expect most of the costs for T1 and T6 to be similar since they fault in many of the same pages. The *min fault* entry in Table 5 is present due to the way our implementation interacts with the virtually mapped CPU cache of the client machine (see Section 3). This effect increased the average fault time by 6% and 5%, respectively for T1 and T6. The entry labeled *page fault* in Table 5 is the time that was required to detect the illegal page access and invoke the fault handler. Page faults comprised 3% of average faulting time for T1 and 2% for T6. We note that it was not possible to measure the times given for the *min fault* and *page fault* entries directly by running the benchmark. These times were obtained instead by measuring a test application that performed the operations several thousand times in a tight loop.

The remaining table entries break down the time spent in the fault handling routine. The entry for *misc. cpu overhead* includes time for looking up the address that caused the fault in the in-memory table, various residency and status checks to determine the appropriate action to take in handling the fault, and other miscellaneous work. *Data I/O* is the time needed to read the page of objects from disk and update the buffer manager's data structures.

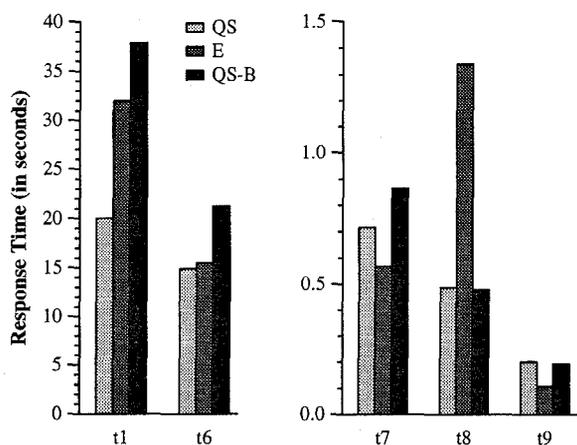


Figure 3. Read-only traversal cold times, small database.

	T1	T6	T7	T8	T9
QS	474	467	26	19	9
E	1018	600	25	18	7
QS-B	1047	639	31	19	9

Table 3. Client I/O requests, small database.

⁴T6: DFS of assembly hierarchy visiting only the root atomic part.

⁵T7: Traverse starting from a randomly selected atomic part up the assembly hierarchy.

⁶T8: Scan the manual object counting occurrences of a specified character.

⁷T9: Compare first and last characters of the manual to see if they are equal.

system	time(ms.)	
	T1	T6
QS	29.4	33.1
E	23.7	26.5
QS-B	31.6	34.5

Table 4. Average Time Per Fault.

description	time(ms.)	
	T1	T6
min faults	1.8	1.6
page fault	.8	.7
misc. cpu overhead	.5	.2
data I/O	24.8	28.5
map I/O	1.1	1.1
swizzling	.4	.4
mmap	.8	.8
total	30.2	33.3

Table 5. Detailed QS Faulting Times.

This accounted for largest fraction of faulting time, 82% for T1 and 85% for T6. The portion of time spent reading mapping tables (*map I/O*) was 3.5% for T1 and 3.2% for T6. The *swizzling* entry gives the time needed to process the mapping table entries. Swizzling costs were quite low, accounting for 1% to 2% of the faulting cost on average. Since all of the pages read were mapped to the locations in memory that they occupied previously, the swizzling time doesn't include any overhead for updating pointers on pages that are inconsistent with the current mapping. The final entry, labeled *mmap*, gives the average time taken by the mmap system call to change the access protections. This accounted for a modest 3% of the faulting time. Finally, we note that the sums of the detailed times given in Table 5 correlate closely with the total per fault times given in Table 4.

We next consider traversals T2 and T3 which include updates. Figure 4 shows the total response time for these traversals run as a

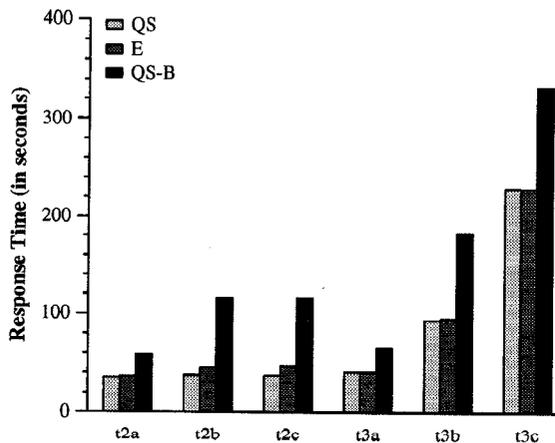


Figure 4. T2, T3 cold times, small database.

single transaction. The I/O requests for T2 were nearly identical to T1 (Table 3), while the T3 traversals performed a few additional I/Os to read index pages. During T2a, which updates the root atomic part of each composite part, QS is 4% faster than E (Figure 4). This may seem surprising given that QS was 37% faster than E during T1 which does the same traversal as T2a, but without updates. The difference in performance between QS and E diminishes during T2a because the page-at-a time scheme for handling updates of QS is more expensive than the object-at-a time approach of E when sparse updates are done.

Part of the increase in response time for QS is due to the fact that the number of page access violations increases from 454 during T1 to 878 during T2a, nearly doubling. The additional access violations occur when the first attempt is made to update an object on a page during the transaction. When this happens, the fault-handling routine is invoked to handle the access violation. As explained in Section 3, this routine performs several main functions. First, it copies the objects contained on the page into a side buffer so that the original values contained in the objects may be used to generate logging information for updates (by diffing) at a later time. Next, it calls ESM, if necessary, to obtain an update (exclusive) lock on the page, and finally, it changes the virtual memory protections on the page, so that the instruction that caused the exception can be restarted. Our measurements showed that during T2a, a total of 5.198 seconds were needed to carry out this work for QS, which amounted to 12.3 ms. for each of the 423 pages updated. Of this, 7.3 ms. was spent copying the objects on the page, 2.8 ms. was used to upgrade the lock on the page, and .9 ms. on average was spent calling *mmap* to change the page's protection to allow write access.

The response time of QS also increases relative to E during T2a because transaction commit is more expensive for QS. The commit time for QS can be broken down into the time required to perform three basic activities, plus a small amount of additional time to perform minor functions like reinitializing data structures, etc. The first of the basic operations involves *diffing* objects on pages that have been updated, and calling ESM to generate log records when it is determined that updates did occur. The diffing phase required a total of 3.035 seconds during T2a, of which .182 seconds was spent calling ESM to generate the 491 log records needed. Thus, the time needed on average to *diff* each of the 423 modified pages (not counting time to generate log records) was 6.7 milliseconds. The second major task performed during transaction commit is updating the mapping tables associated with each modified page. Our measurements showed that 3.084 seconds (7.2 ms. per page) were required for this phase of commit processing. The final step in committing a transaction is performed by ESM. This involves writing all log records to disk at the server, and flushing all dirty pages back to the server from the client. This phase of commit processing required 3.501 seconds during T2a.

Turning to T2b and T2c, we see that QS is 17% and 20% faster than E, respectively. As one would expect, QS does better relative to E during T2b and T2c when updates are more dense since QS copies and *diffs* fewer objects unnecessarily. In fact the performance of QS degrades only slightly during T2b relative to T2a. This is due almost entirely to increased time during commit for *diffing* objects and generating log records. More precisely, during T2b 5.804 seconds was required do the *diffing* (.280 seconds of this was for generating log records). The average *diffing* cost per page was 12.9 ms during T2b (not counting logging). We also note that the performance of QS was basically the same during T2b and T2c, while the performance of E was 5% slower. This is because repeatedly updating an object is very cheap for QS since objects are accessed via virtual memory pointers while updating an object in E requires a function call.

The performance difference between QS and E narrows further during T3 relative to T2 and T1. QS has better performance than E in all cases, but nearly similar overheads for index maintenance make this difference less noticeable. In contrast to the relatively stable performance of the systems during T2, the response times of the systems steadily increase when going from T3a to T3b to T3c. This is because each update of an indexed attribute results in the immediate update and logging of the update to the corresponding index. QS-B is always much slower than the other systems during T2 and T3, especially during the B and C traversals. This is because the 4Mg area used to hold recovery data wasn't big enough to hold all of the objects from modified pages during these traversals for QS-B which caused additional log records to be generated.

5.3. Small Hot Results

The hot results were obtained by re-running the OO7 benchmark operations after all of the data needed by each operation had been cached in the client's memory by the cold traversal. Figure 5 shows hot times for the traversals run on the small database. The times for QS-B are omitted since they are identical to those shown for QS.

As one would expect, the performance of QS is generally better than E. It is somewhat surprising, however, that E is just 23% slower than QS during T1. To determine the reasons for this relatively small difference, we used *gpt*[Ball92] to profile the benchmark application. Table 6 presents the results of the profiling for T1. The T1 hot traversal time has been broken down in Table 6 based on the percentage of CPU time spent in several groups of

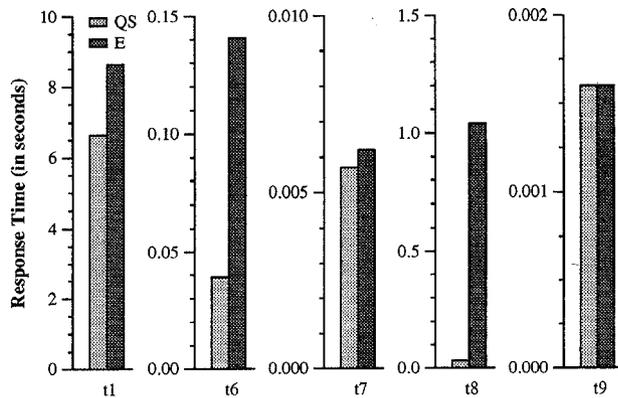


Figure 5. Traversal Hot Times.

description	% of time	
	QS	E
EPVM 3.0	-	33.31
malloc	56.13	24.99
part set	35.18	24.57
traverse	8.03	17.12
do nothing	0.64	0.70
misc.	0.02	0.01
total	100.00	100.00

Table 6. T1 hot traversal detail.

functions. Table 6 shows that E spent 33% of the time executing EPVM 3.0 interpreter functions. Most of this time was spent dereferencing unswizzled pointers. Both QS and E spent a considerable amount of time allocating and deallocating space in the transient heap (see the entry for *malloc*). This is because an "iterator" object is allocated in the heap for every node (assembly object, composite part, and atomic part) in the object graph that is visited during the traversal. The "iterator" object establishes a cursor over the collection of pointers to sub-objects so that the sub-objects can be traversed.

The entry labeled *part set* in Table 6 gives the time spent in functions that maintain a set of the atomic part ids visited in each composite part's subgraph of atomic parts. This set is needed so that the same atomic part is not visited more than once by the DFS traversal. The time spent in other functions that implement the traversal, such as functions that iterate over collections of pointers to sub-objects and that implement the recursive traversal was 8% for QS and 17% for E. The higher percentage for E reflects the additional cost of dereferencing large pointers in E. When each node in the object graph is visited, a simple function is called that examines a field in the object to make sure that the object is faulted into memory. The time spent in these functions was .7% for both systems. The detailed numbers in Table 6 are surprising because the small amount of additional work involving transient data structures that was needed to implement T1 accounts for a such a large percentage of the overall cost. The results show how quickly a small amount of additional computation can mask differences in the cost of accessing persistent data between the systems.

E is 3.6 times slower than QS during T6. QS performs better relative to E during T6 (the sparse traversal) than during T1 (the dense traversal) since there is less overhead for maintaining transient data structures during T6. For example, the sets of part ids are not maintained since only the root part is visited during T6. The performance of the systems is very close during T7. T7 visits very few objects in the database (10 to be precise) since it simply follows pointers from a single atomic part up to the root of the module. Thus, differences in traversal cost are easily diminished by other costs, such as the overhead for looking up the atomic part up in the index to begin the traversal, etc. Figure 4 shows that E is a factor of 32 slower than QS during T8. T8 scans the manual object, a large object spanning several pages on disk. In the case of E, an EPVM 3.0 function call is performed for each character of the manual that is scanned while QS accessed each character of the manual via a virtual memory pointer. Profiling showed that E spent 91% of its time executing EPVM functions during T8. During T9 the systems have identical performance. Profiling showed that the hot results for T9 largely reflect similarities in things such as index lookup costs between the systems since an index lookup is performed to locate the module object.

5.4. Medium Cold Results

This section presents the cold times for the OO7 benchmark operations when run on the medium database. The results presented represent the average of 5 runs of the benchmark experiments. Figure 6 presents the cold response times for the traversal operations and Table 7 gives the number of client I/O requests. We see in Figure 6 that, as in the case of the small database, QS has the best performance during T1. QS is 41% faster than E during T1 while it performs 63% fewer I/Os. E, on the other hand, is 36% faster than QS-B during T1.

E has better performance than QS during T6 and T7. QS is slower during T6 because the number of page-faults between the two systems is similar and QS has higher costs per fault. The relative times shown for T7 and T8 are close to the small database case. QS is slower due to higher per fault costs during T7. E is slower than QS during T8 due to the overhead of calling EPVM to scan

each character of the manual. The results for T9 (not shown) were identical to the small case.

Turning now to traversals T2 and T3 (Figure 7) which perform updates, we see that QS outperforms E during the T2a and T3a traversals which only update the root atomic part of each composite part. This is understandable when one considers that both QS and do basically the same amount of work to process the updates that they did in the small case for T2a and T3a. This makes the cost difference for doing the traversal itself the main factor effecting their relative performance. The relative performance of QS worsens during T2b and T2c causing QS and E to have similar performance. Recovery is more expensive for QS during T2b and T2c since the buffer used for recovery is much smaller than the fraction of the database that is updated. QS-B has much worse performance than both QS and E in Figure 7. This is caused by the fact

that in addition to higher traversal costs, QS-B has higher costs for recovery as well.

5.5. Effect of Relocating Pages

Recall that QuickStore always tries to assign a disk page to the virtual memory locations that it last occupied when in memory. This section considers the effect on performance of relocating pages at different memory addresses. This increases faulting costs because pointers between persistent objects then have to be updated to reflect the new assignment of disk pages to virtual memory addresses. We consider two approaches to dealing with page relocations. The first approach updates or swizzles pointers that need to be modified when pages are faulted into memory, but these changes are not written back to the database. This implies that the changes will have to be made again if the same data is accessed in subsequent program runs. We refer to this system as QS-NW. The second approach commits the changed mapping to the database. This approach is more costly initially, but may be able to avoid further relocations in the future. This approach also has the disadvantage that it can turn a read-only transaction into an update transaction. We refer to this approach as QS-WR.

Table 8 presents the results for T1 run on the small database when the percentage of pages that are relocated in memory is varied from 0 to 100%. The pages that were relocated in the experiment were picked at random. Table 8 shows that when the number of relocations is small (5%), the performance of the systems is not greatly affected, however, when the relocation percentage is 20%, QS-WR is 25% slower than when no relocations occur. The difference in performance between the two schemes at this point is also about 25%. The performance of QS-NW slows by 7% and 38% when the percentage of relocated pages is 50% and 100% respectively, while QS-WR suffers a 67% reduction in performance when the relocation probability is 50%. QS-WR is much slower than QS-NW when all pages are relocated since it must commit updates for all of the pages in the database.

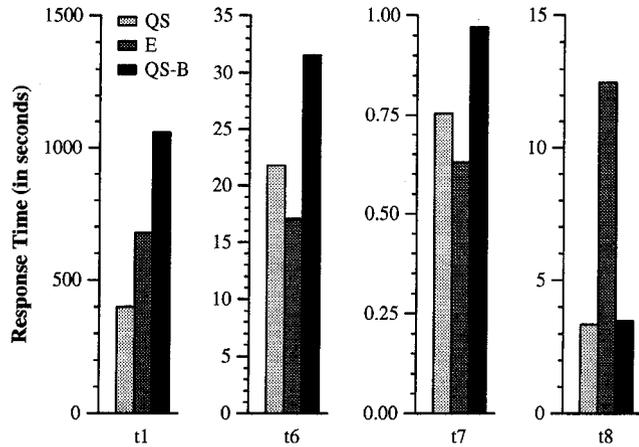


Figure 6. Cold Times, Medium Database.

	T1	T6	T7	T8
QS	13216	610	27	130
E	35622	558	25	129
QS-B	36963	802	32	130

Table 7. Traversal Cold I/Os.

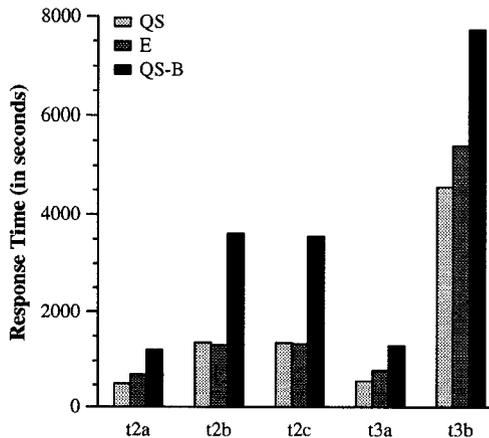


Figure 7. T2, T3 cold times, medium database.

6. Conclusions

This paper has presented the design of QuickStore, a memory-mapped storage system built on top of the Exodus Storage Manager. The paper also compared the performance of QuickStore with E, a persistent version of C++ that uses a software interpreter to support access to persistent data. The OO7 benchmark was used as a basis for comparing the performance of the two systems. The purpose of the study was to accurately measure the differences in performance of the hardware and software based pointer swizzling schemes. The results of the study give a clear picture of the tradeoffs between the two approaches, which we summarize below.

The results of the cold traversal experiments showed that when object accesses are dense (T1), QuickStore has the best performance. This is because object sizes in QuickStore are smaller than in E, due to the different schemes used by the systems to represent pointers on disk. Its smaller object size allowed QuickStore to perform significantly fewer disk I/O operations to read the same amount of data when accesses were clustered. When object

	0%	5%	20%	50%	100%
QS-NW	20.085	19.748	20.252	21.454	27.842
QS-WR	20.085	20.183	25.066	33.594	60.297

Table 8. T1 response time (in seconds), vary % of relocations.

accesses were unclustered (T6), the performance of QuickStore was comparable (small database) or worse (medium database) than the performance of E. The reason for this change, relative to the clustered case, was that there was less difference between the systems in the number of pages faulted into memory per object access during the unclustered experiments. This exposed the fact that QuickStore has higher per faults costs than E. In addition, there were some cases (T7 and T9) when QuickStore always had slower performance due to higher faulting costs. The slower performance for QuickStore during T7 was influenced by the cost of reading a relatively large amount of mapping information to support the memory-mapping scheme that it uses.

The higher faulting costs for the memory-mapping scheme were also highlighted by the performance of QS-B (QuickStore with big objects). QS-B always had slower performance than E during the read-only cold experiments, except during the experiments that manipulated large objects (T8). The memory-mapped schemes had better performance than E when large objects were accessed because large object accesses require significantly more CPU work under the software approach. This additional cost caused E to be slower even in the cold case.

For the traversals in which faulting costs were examined in detail, it was shown that the average cost per fault for QuickStore was roughly 20% higher than for E. The largest component of the additional faulting cost for the memory mapping scheme was the time required to read mapping information from disk. This comprised 4% of the average cost per fault. The detailed cost analysis also showed that the overhead for handling page protection faults and manipulating page access protections were each 3%. The smallest component of the faulting cost for QuickStore was the CPU cost for swizzling pointers. This was just 1% of the average cost per fault.

The performance of QuickStore was generally better than E when updates were performed. The results of the update experiments showed, however, that the page-based diffing scheme used by QuickStore to generate log records was more expensive when updates were sparse and when the update activity was heavy enough to cause log records to be generated before transaction commit. QuickStore performed better relative to E when a higher percentage of objects were updated on each page since QuickStore copied and diffed fewer objects unnecessarily in this case. The detailed times for the update experiments showed that the cost of diffing was ranged from 7 to 12 milliseconds per page.

The hot results helped to quantify the performance advantage of the memory-mapped scheme when working on in-memory objects. In some cases (T1) the difference in performance between QuickStore and E was only 23%, while in others (T6) QuickStore was over 3 times faster than E. This showed how quickly the performance of the systems converged when a small amount of additional work was performed. The results also showed that E was significantly slower than QuickStore when doing in-memory work on large objects since this required all accesses to be handled by the E interpreter.

We also examined the performance of QuickStore when pages of objects must be relocated in memory. This increases the amount of swizzling work performed by QuickStore. When the percentage of pages that were relocated was small, the performance of the systems did not noticeably worsen. However, a high percentage of relocations did have a noticeable effect on overall performance. In particular, when the new mapping tables were written back to the database, performance worsened by a factor of three.

In the future we would like to consider the performance impact of different swizzling approaches on query workloads. Queries tend to have sparse access patterns, so systems that do pure hardware swizzling may not perform well during queries. We are also

interested in the impact of versioned data on the performance of different pointer swizzling techniques. It would also be interesting to investigate alternatives to the page-based diffing approach used by QuickStore to support recovery. For example, the approach used by ObjectStore is to log entire pages of modified objects. The performance of the diffing approach could also be made more efficient if some level of compiler support were available.

References

- [Ball92] T. Ball and J. Larus, "Optimally Profiling and Tracing Programs", POPL 1992, pp. 59-70, January 1992.
- [Carey89a] M. Carey et al., "The EXODUS Extensible DBMS Project: An Overview," in *Readings in Object-Oriented Databases*, S. Zdonik and D. Maier, eds., Morgan-Kaufman, 1989.
- [Carey89b] M. Carey et al., "Storage Management for Objects in EXODUS," in *Object-Oriented Concepts, Databases, and Applications*, W. Kim and F. Lochovsky, eds., Addison-Wesley, 1989.
- [Exodu92] Using the EXODUS Storage Manager V2.0.2, technical documentation, Department of Computer Sciences, University of Wisconsin-Madison, January 1992.
- [Carey93] M. Carey, D. DeWitt, J. Naughton, "The OO7 Benchmark", Proc. ACM SIGMOD Int'l Conf. on Management of Data, Washington, DC, May 1993.
- [Frank92] M. Franklin et al., "Crash Recovery in Client-Server EXODUS", Proc. ACM SIGMOD Int'l Conf. on Management of Data, San Diego, California, 1992.
- [Hoski93] A. Hosking, J. E. B. Moss, "Object Fault Handling for Persistent Programming Languages: A Performance Evaluation", OOPSLA '93, pp. 288-303
- [Lamb91] C. Lamb et al., "The ObjectStore Database System", CACM, Vol. 34, No. 10, October 1991
- [Moss90] J. Eliot B. Moss, "Working with Persistent Objects: To Swizzle or Not to Swizzle", COINS Object-Oriented Systems Laboratory Technical Report 90-38, University of Massachusetts at Amherst, May 1990.
- [Objec90] Object Design, Inc., ObjectStore User Guide, Release 1.0, October 1990.
- [Rich93] J. Richardson, M. Carey, and D. Schuh, "The Design of the E Programming Language", ACM Trans. on Programming Languages and Systems, Vol. 15, No. 3, July 1993.
- [Rich90] J. Richardson, "Compiled Item Faulting", *Proc. of the 4th Int'l. Workshop on Persistent Object Systems*, Martha's Vineyard, MA, September 1990.
- [Schuh90] D. Schuh, M. Carey, and D. DeWitt, "Persistence in E Revisited--Implementation Experiences, in *Implementing Persistent Object Bases Principles and Practice*", The 4th Int'l. Workshop on Persistent Object Systems.
- [Shek90] E. Shekita and M. Zwilling, "Cricket: A Mapped Persistent Object Store", *Proc. of the 4th Int'l. Workshop on Persistent Object Systems*, Martha's Vineyard, MA, Sept. 1990.
- [Singh92] V. Singhal, S. Kakkad, and P. Wilson, "Texas: An Efficient, Portable Persistent Store", in *Proc. of the 5th Int'l. Workshop on Persistent Object Systems*, San Miniato, Italy, Sept. 1992.
- [White92] S. White and D. DeWitt, "A Performance Study of Alternative Object Faulting and Pointer Swizzling Strategies", in *Proc. of the 18th Int'l. Conf. on Very Large Data Bases*, Vancouver, British Columbia, August 1992.
- [Wilso90] Paul R. Wilson, "Pointer Swizzling at Page Fault Time: Efficiently Supporting Huge Address Spaces on Standard Hardware", Technical Report UIC-EECS-90-6, University of Illinois at Chicago, December 1990.