

# An Adaptable Multithreaded Prefetching Technique for Client-Server Object Bases

Nils Knaffa

*Department of Computer Science  
The University of Edinburgh  
James Clerk Maxwell Building  
King's Buildings  
Mayfield Road  
Edinburgh EH9 3JZ, Scotland  
E-mail: nk@dcs.ed.ac.uk*

Given the existence of powerful multiprocessor client workstations in many client-server object database applications, the performance bottleneck is the delay in transferring pages from the server to the client. We present a prefetching technique that can avoid this delay, especially where the client application requests pages from several database servers. This technique has been added to the EXODUS storage manager. Part of the novelty of this approach lies in the way that multithreading on the client workstation is exploited, in particular for activities such as prefetching and flushing dirty pages to the server. Using our own complex object benchmark, we analyze the performance of the prefetching technique with multiple clients and multiple servers. The technique is also tested under a variety of client host workload levels.

**Keywords:** prefetching, object-oriented databases, distribution, performance analysis, multithreading, application access pattern, storage management

## 1. Introduction

Much of the research and development effort in high-performance database systems has focused on exploiting parallel computing on the database server platform. However, with the rapidly improving performance/price ratios for shared-memory multiprocessors it is now feasible to use parallel computers as client hosts. This raises the question of how a multiprocessor client machine can be exploited to improve database response time. In one approach [8] used so-called

*ParSets* in the OO7 benchmark traversals [4] to invoke a method on every object in a set in parallel. This approach can boost performance where the application is CPU-bound. However, for some applications, the performance bottleneck is dominated by the delay that results from the client waiting for the transfer of pages from the server and writing pages to the server.

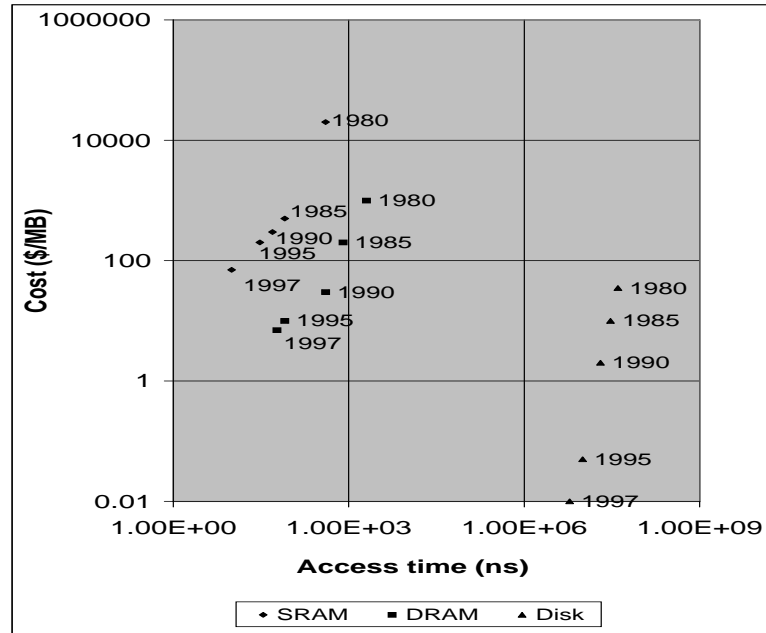


Figure 1. Performance/price development

In this paper, we examine the way in which prefetching can exploit parallelism. Prefetching has been studied before in many areas of computing such as operating systems, microprocessor design, compiler construction, the world wide web and databases. Although prefetching has been studied for a long time, the problem prefetching is facing has changed over the years as a result of technology advances. Fig. 1 shows the performance/price development of semiconductor memories and magnetic disks. There is a two-order of magnitude gap in access time between memory and disks. Memory access is faster and the rate of improvement is also higher. The prices of memory are falling which makes database

caches cheaper and buffer replacement less of a problem. Prices of disks are also falling dramatically which is good for cheap secondary storage, allowing increased use of RAIDs.

CPU performance improves at an even higher rate than memory (fig. 2). CPU is doubling its performance every 18 months whereas disk retrieval time only improves about 5% per year. Memory access time improves at about 10 to 12% per year. Client workstations will become powerful multi-processors machines with high speed CPUs, most of which will tend to be idle most of the time.

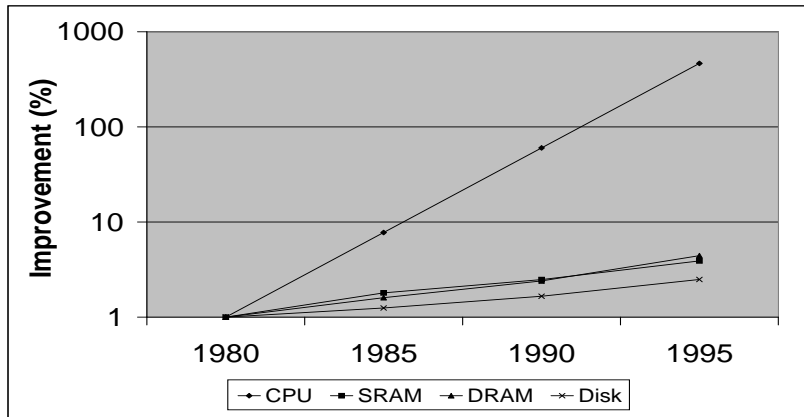


Figure 2. Performance improvement in %

Gray [12] asked “What happens when processors are infinitely fast and storage is free?” He predicts that technology trends promise to give us processors with pico-second clock speeds. These pico-processors will spend much of their time waiting for information from the storage hierarchy. Similar trends will bless us with peta-byte online stores with exa-byte near-line stores.

So the client-server system of the very near future is one in which the client is a chronically under-utilised multiprocessor that spends most of its life waiting for disk blocks to arrive from a remote server.

We address this problem by using multithreading to parallelise the I/O for prefetching and flushing. Various techniques of buffer management [1], [2] have been used to reduce traffic between the client and the server. However, in this

paper, we concentrate on techniques that attempt to effect page transfers in a more timely fashion.

The fate of object-oriented database management systems (OODBMSs) will largely depend on their performance in comparison to relational databases. The simple tabular structures of relational databases and the set-at-a-time semantics of retrieval languages such as SQL make it easy to parallelise relational database servers. However, in an OODBMS the structures are complex and typically the retrieval chases pointers. Furthermore, in most OODBMSs the bulk of the processing occurs on the client: the server merely serves pages. With the increasingly powerful client workstations OODBMSs take a lot of work from the server to the client.

Early work on prefetching in relational databases [18] was able to exploit sequential access patterns. In object-oriented databases, the object relationships are often used to predict future accesses. The object structure and the pointer relationships provide important information for object access patterns.

Chang and Katz [5] predict future accesses from the data semantics in terms of inheritance and structural relationships. Objects are stored in pages and they prefetch the immediate object reference only. Due to the high cost of a page fetch, this technique has only a limited effect on reducing elapsed time. Cheng and Hurson [6] extended this work by adding multiple hints, a prefetching depth and physical storage considerations. Instead of using a single hint, a series of hints are given for all types of relationships. Prefetch depths were added to the prefetching hints according to the semantics of each relationship. For example, an application may require the access to follow configuration links recursively or to follow version links at a maximum of three levels away. Physical storage considerations are used to impose a limit on high cost I/O.

A complex assembly operator to load component objects recursively in advance was introduced by Keller et al [14]. The application traversal was performed by three different scheduling algorithms, *depth-first*, *breadth-first* and an *elevator* algorithm (which schedules disk access to objects based on their physical location). Objects were clustered (according to their type or to the composite object structure) or unclustered. This technique is suitable for small object nets; in larger object nets, recursively prefetched objects might already have been replaced by the buffer replacement strategy.

In the Fido system [17], the prefetching technique employs an associative memory to recognize access patterns within a context over time. In training

mode, object access information is gathered and stored with a *nearest-neighbor* associative memory. In prediction mode, this information is used to recognize previously encountered situations.

Gerlhof developed an architecture for prefetching [10] and a so-called *Prefetch Support Relation* (PSR) [11]. The PSR stores the precomputed page answer of an operation, i.e. the identifiers of all pages that were accessed during the execution of an operation. [13] developed an extensible file system, ELFS, in which they used user hints to predict the file access pattern.

In Thor [7] each fetch request from the client causes the server to select a prefetch group containing the object requested and possibly some other objects. A fetch request is processed to completion, determining all members of the prefetch group, before any objects are sent to the client. The disadvantage of this approach is that the server is already a bottleneck and additional prefetch requests further increase the server workload.

Our prefetching technique observes the client processing on the *object net*. If the application moves towards a non-resident object, the page of the object is a candidate for prefetching. We try to time the prefetch request so that the request is not started too early but the page arrives before application access. We implemented this technique by extending the EXODUS storage manager (ESM) [3].

In section 2 we describe the design of our prefetching technique and of our buffer replacement strategy. An overview of ESM and the prefetching architecture is given in section 3. In section 4 we present the theoretical results and present the empirical results from our benchmark. Finally, in section 5 we conclude our work and give an outlook for future work.

## 2. The Prefetching Design

### 2.1. Prefetch Object Table

OODBMSs can store and retrieve large, complex data structures which are nested and heavily interrelated. Examples of OODBMS applications are CAD, CAM, CASE and Office automation. These applications consist of objects and relationships between objects containing a large amount of data. A typical scenario is laid out by the OO7 benchmark [4]. It comprises a very complex assembly object hierarchy and is designed to compare the performance of object-oriented

databases.

In a page server, like ESM, objects are clustered into pages. Good clustering is achieved when references to objects in the same page are maximized and references to objects on other pages are minimized. In our benchmark we use a *composite object clustering* technique.

The general idea of our technique is to prefetch references to other pages in a complex object structure net (e.g. OO7). We obtain the prefetch information from the object references without knowledge of the object semantics. Considering the object structure in a page, we identify the objects which have references to other pages (*Out-Refs*). One page could possibly have many *Out-Refs* but sometimes it is not possible to prefetch all pages because of time and resource limitations. Instead, we observe the client navigation through the object net. We know which objects have *Out-Refs* and when we identify that the application is processing towards such an *Out-Ref-Object* (ORO) the *Out-Ref* page becomes a candidate for prefetching.

The prefetch starts when the application encounters a so-called *Prefetch Start Object* (PSO). Although the determination of OROs is easy, determining PSOs is slightly more complicated. There are two factors that complicate finding PSOs:

1. Prefetch Object Distance (POD)

For prefetching a page it is important that the prefetch request arrives at the client before application access to achieve a maximum saving. The POD defines the optimal distance of  $n$  objects from the PSO to the ORO object which is necessary to provide enough processing to overlap with prefetching. Let  $C_{pf}$  denote the cost of a page fetch and let  $C_{op}$  denote the cost of object preparation. The cost of object preparation is the ESM client processing time before the application can work on the object<sup>1</sup>. Then POD is computed as follows:

$$POD = \frac{C_{pf}}{C_{op}}$$

If the prefetch starts before the POD, a maximum saving is guaranteed, however, if it starts after the POD, but before access, some saving can still be achieved (see section 4.2).

<sup>1</sup> Additionally we could use the expected amount of processing from the application.

## 2. Branch Objects

A complex object has references to other objects. The user of the application decides at a higher level the sequence of references with which to navigate through the object net. We define a *Branch Object* as an object which has at least two references to other objects. Objects that are referenced by a *Branch Object* are defined as a *Post-Branch Object*. For example in fig. 3 we have a complex object hierarchy. The object with the OID (Object Identifier) 1 would be defined as a *Branch Object* because it contains a branch in the tree of objects. Objects with OID 2, OID 7 and OID 12 would be defined as *Post-Branch Objects* because they are the first objects de-referenced by a *Branch Object*.

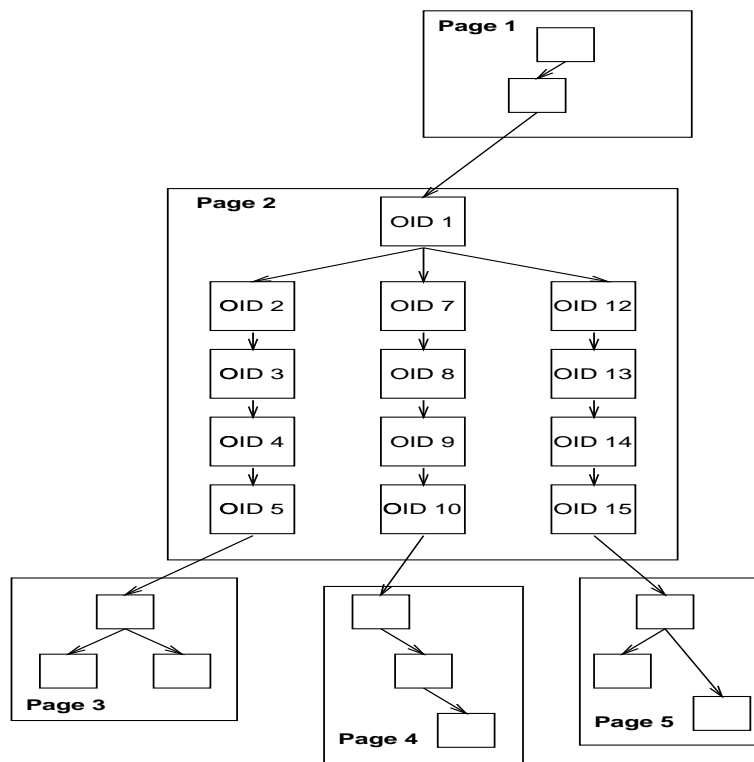


Figure 3. Object relationship

For every identified ORO in the page we compute the PSO by the following algorithm:

1. Retrieve the OID of the ORO and of the object in the next page referenced by the ORO.
2. Compute the POD to define the distance of  $n$  objects from PSO to ORO.
3. Determine the PSO by following the object reference  $n$  objects backwards from the ORO. If there are not enough objects in the reference chain before the ORO, then we will identify the first object in the page from the reference chain to achieve at least some saving.
4. If the object is already identified as a PSO and the previously identified PSO has a different *Post-Branch Object* then we would identify the *Post-Branch Objects* of the object as PSOs<sup>2</sup>.

Defining *Post-Branch Objects* as PSOs can improve the accuracy for the prediction and reduces the number of adjacent pages to prefetch. For example in fig. 3 we would identify OID 5, OID 10 and OID 15 as OROs. In this example we assume a POD of 4 objects. On analyzing page 2 we would identify OID 5 as an ORO. From OID 5 we would go through the chain backwards by 4 objects and identify OID 1 as a PSO. Then we would do the same for the OROs OID 10 and OID 15 and identify OID 1 as the PSO for both. After analyzing the whole page we would find out that OID 1 has three PSOs with different *Post-Branch Objects*. In this case we would identify the *Post-Branch Objects* of OID 1 (OID 2, OID 7 and OID 12) as PSOs instead of OID 1.

The novel idea about our technique is to make prefetching adaptable to the client processing on the object net. Because the cost of a page fetch is high we try to start the prefetch early enough to achieve a high saving but not too early to prefetch inaccurately. In contrast to the work of [14], we do not prefetch all references recursively; instead we select the pages to prefetch, dependent on the client processing. Recursive object prefetching has also the problem that prefetched pages can be replaced again before access. Adaptive object prefetching limits the number of prefetch pages to the adjacent pages. In contrast to [5], we look further ahead for objects to prefetch than the immediate object.

Each page of the database is analyzed off-line. The Analyzer stores this information in the POT for every database root<sup>3</sup>. Fig. 4 depicts the layout of

<sup>2</sup> This step is executed after we have defined all PSOs from the OROs in a page.

<sup>3</sup> This is important because objects on the same page could belong to different roots.



<b>Page</b> [PageID]	<b>PSO</b> [SlotIndex]	<b>ORO</b> [SlotIndex]	<b>RefPage</b> [PageID]	<b>RefOID</b> [SlotIndex]
-------------------------	---------------------------	---------------------------	----------------------------	------------------------------

Figure 4. Entry in POT

one entry in the POT. Entries for one page are clustered together on disk. The overhead for this table is quite low as it only contains a few objects of the page.

At run time, the information from the POT is used to start the prefetch requests. The run time system allocates enough threads for prefetching. If essential pointers for the navigation are updated in a transaction we would invalidate the POT for this page and modify it after the completion of the transaction.

This prefetching technique is not only useful for complex objects, it can also be used for collection classes (linked list, bag, set or array) in OODBMSs. Applications traverse an object collection with a cursor. With PSO and ORO it would be possible to prefetch the next page from a cursor position. In the description of our technique, the object size is assumed to be smaller than the page size. If the object is larger than a page, prefetching can be used to bring the whole object into memory.

In future work we want to investigate performance and behaviour of the object navigation using probability values. If the POD has many objects and the fan-out of the objects is high then the probability that the navigation will be from the PSO to the ORO is low. We will compare performance under a variety of assumptions (a) the probability of the navigating to the ORO and (b) the benefit of a correct prefetching and the cost of an incorrect prefetch.

## 2.2. Replacement Policy

In the ESM client it is possible to open buffer groups with different replacement policies (LRU and MRU). Freedman and DeWitt [9] proposed a LRU replacement strategy with one chain for demand reads and one chain for prefetching. We also plan to use two chains with the difference that when a page in the demand chain is moved to the top of the chain, the prefetched pages for this page are also moved to the top. The idea of this algorithm is that when the demand page is accessed, it is likely that the prefetched pages are accessed too. If a page from the prefetch chain is requested it is moved into the demand chain.

### 3. System Architecture

#### 3.1. The EXODUS Storage Manager

We implemented the prefetching technique in ESM. The EXODUS client-server database system [3] was developed at the University of Wisconsin. It aids a database implementor in the task of generating a DBMS by providing a storage manager, a programming language E (an extension of C++), a library of access-method implementations, a rule-based query optimizer generator, and tools for constructing query-language optimizers.

The basic representation for data in the storage manager is a variable-length byte sequence of arbitrary size, incorporating the capability to insert or delete bytes in the middle of the sequence. In the simplest case, these basic storage objects are implemented as contiguous sequences of bytes. As the objects become large, or when they are broken into non-contiguous sequences by editing operations, they are represented using a B-tree of leaf blocks, each containing a portion of the sequence. Objects are referenced using structured OIDs<sup>4</sup>.

On these basic storage objects, the storage manager performs buffer management (LRU or MRU), concurrency control, recovery, and a versioning mechanism that can be used to provide a variety of application-specific versioning schemes. Transactions are implemented using a shadowing and logging technique. Client and server communicate via the socket interface. The client specifies the requested data in a message structure and sends it to the server. The server updates this structure and responds with the attached 8K page.

#### 3.2. The Prefetching Architecture

In this section we describe how prefetching is incorporated into ESM. For the concurrent execution of the application and the prefetch system we used the Solaris thread interface [19]. Multithreading combined with prefetching has the benefits of:

1. Increased application throughput and responsiveness;
2. Performance gains from multiprocessing hardware (parallelism);
3. Efficient use of system resources.

<sup>4</sup> Object identifiers containing a physical and logical component (in ESM page number and slot number).

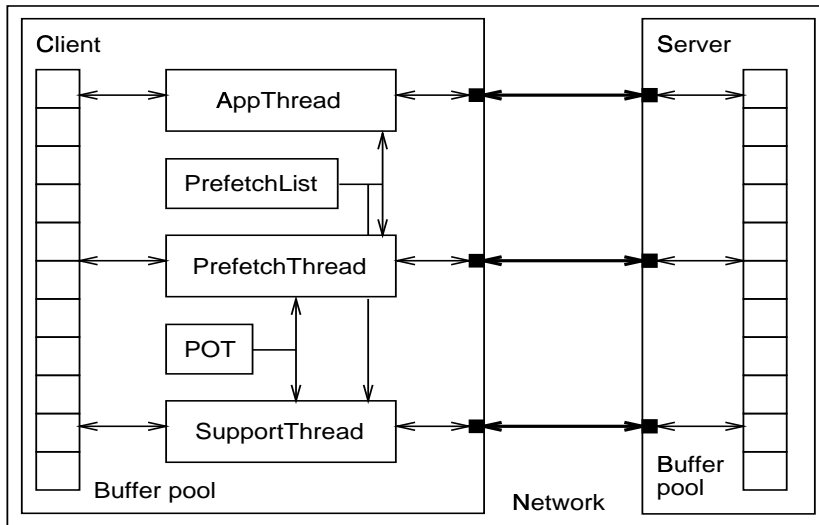


Figure 5. Prefetching architecture

As depicted in fig. 5, the database client is multithreaded. The *AppThread* is responsible for the processing of the application program and the *PrefetchThread* is responsible for fetching pages in advance into the buffer pool. A *SupportThread* has the same task as the *PrefetchThread* with the only difference being that it is scheduled by the *PrefetchThread*. Each thread has one associated socket. The POT informs the *PrefetchThread* which pages are candidates for prefetching from the current processing of the application. The *Prefetch List* is a list of pages which are currently prefetched. A *FlushThread* is responsible for flushing dirty pages to the server.

At the beginning of a transaction the *AppThread* requests the first page from the server by a demand read. The *PrefetchThread* always checks which objects the *AppThread* is processing. Having obtained this information, it consults the POT for a page to prefetch and checks if this page is already resident. If not, the page is inserted in the *Prefetch List* and the request is sent to the server. The server responds with the demanded page and the client inserts the page into its buffer pool. Eventually the page is removed from the *Prefetch List* and inserted into the hash table of the buffer pool.

If the POT predicts multiple pages, *SupportThreads* help the *PrefetchThread*; this is useful when the prefetch object distance is short. The number of *SupportThreads* is determined by the number of simultaneous prefetch requests. Each

*SupportThread* runs on its own LWP<sup>5</sup> and while one *SupportThread* blocks on I/O, another *SupportThread* can insert its page into the buffer pool.

When the *AppThread* requests a new page, it first checks if the page is in the buffer pool. If the page is not resident then it checks the *Prefetch List*. If the page has been prefetched the *AppThread* waits on a semaphore until the page arrives, otherwise it sends a demand request to the server.

The ESM server is not multithreaded<sup>6</sup> and performs each request sequentially. But the server forks a new process for the disk management. The server and disk manager communicate via shared memory. The server puts a request for a new page in a disk queue and the disk manager reads the page from disk and copies it into the buffer pool of the server. Incorporating threads into the server would further improve the whole systems performance and is part of future work.

For the parallel execution of threads on the client, synchronization mechanisms are required. The access to the buffer pool is protected by mutexes, which means that only one thread at a time is able to make a residency check or manipulation. When either the *AppThread* or *PrefetchThread* are idle they wait on a semaphore. A semaphore informs the *SupportThread* that there is a page to prefetch.

Prefetch threads are mostly idle as they await the completion of I/O. This means that several threads can be allocated to a single processor and the threads will not have to wait for an operating system time-slice to complete before they can execute. The Solaris thread interface provides a function to give the threads priorities. The *AppThread* has the highest priority to make sure that the application processing always gets scheduling priority on one of the CPUs before the prefetch threads.<sup>7</sup> The *PrefetchThread* has a 50 percent priority<sup>8</sup> and the *SupportThreads* have low priorities.

<sup>5</sup> Lightweight process (LWP) can be thought of as a virtual CPU that is available for executing code.

<sup>6</sup> But ESM runs many tasks, as concurrent processes, on one processor.

<sup>7</sup> On a uniprocessor, a subtler approach to allocating priorities would be needed in order to strike a balance between application processing and prefetching.

<sup>8</sup> Priorities in Solaris are integer values from 0 to 127.

Table 1  
Computer performance specification

Parameter	Server	Client
SPARCstation	20 Model 612	10 Model 514
Main Memory	192 MB	224 MB
Virtual Memory	624 MB	515 MB
Number of CPUs	2	4
Cycle speed	60 MHz	50 MHz

Table 2  
Disk controller performance

Parameter	Disk controller
External Transfer Rate	9 Mbytes/sec
Average Seek (Read/Write)	8 msec
Average Latency	4.17 msec

## 4. Performance Evaluation

### 4.1. System Environment

In table 1 we give a specification of the computers used in our experiments. The client machine has 4 processors. The Ethernet network is running at 10 Mb/sec. The performance of the disk controller (Seagate ST15150W) is presented in table 2.

### 4.2. Theoretical Results

The success of prefetching is dependent on the accuracy of the prediction and the completion of the prefetch before access. We define the cost of object processing to be  $C_o$ . Let  $C_{op}$  denote the cost of preparing one object for application access and let  $C_{oa}$  denote the cost of processing on the object from the application plus waiting time.  $C_o$  is calculated by:

$$C_o = C_{op} + C_{oa} \quad (1)$$

The cost of a page fetch,  $C_p$ , is dependent on client and server processing and the network.  $Cl_b$  denotes the cost of client processing on the buffer pool;  $N_t$  denotes the cost of network transfer;  $S_b$  is the cost of server processing on the

buffer pool;  $S_q$  is the server queueing cost and  $S_d$  is cost for the disk access.  $C_p$  is calculated by:

$$C_p = Cl_b + N_t + S_b + S_q + S_d + N_t + Cl_b \quad (2)$$

The saving for one out-going reference  $S_{or}$  is dependent on the number of objects between the start of the prefetch and application access to the prefetched object ( $N_o$ ) and  $C_p$ :

$$S_{or} = \begin{cases} C_p & \text{if } (C_o \cdot N_o \geq C_p) \\ C_o \cdot N_o & \text{otherwise} \end{cases} \quad (3)$$

If there is enough processing (i.e.  $C_o \cdot N_o$ ) to overlap then the saving is the cost of a page fetch. If not, there is also, albeit lower, saving of the amount of processing from prefetch start to access ( $C_o \cdot N_o$ ). Pages normally have many out-going references. The number of references to different pages is denoted by  $n$ .  $S_p$ , the saving for a whole page, is given by:

$$S_p = \sum_{i=1}^n S_{or}(i) \quad (4)$$

Finally, the saving of the total run is defined by  $S_r$  which is influenced by the cost of the thread management ( $C_t$ ), by the cost of the socket management ( $C_s$ ) and by the number of pages in the run ( $q$ ):

$$S_r = \left( \sum_{j=1}^q S_p(j) \right) - C_t - C_s \quad (5)$$

### 4.3. Benchmark description

For the evaluation of the prefetching technique we created a benchmark with complex objects. The requirements for the benchmark were:

- The application access pattern should be dynamic and different for every run;
- The sizes of the objects should be fairly uniform;
- Object references should be complex;

- The number of pages accessed in one run should be equal to, or less than, the number of pages in the buffer pool at the client and server.

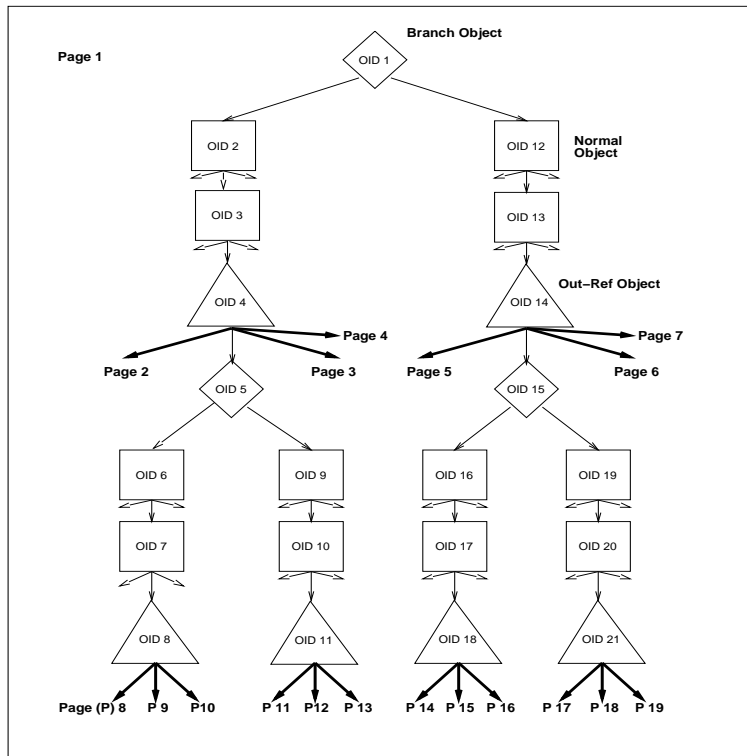


Figure 6. Benchmark structure of one page

In fig. 6 we depict the design of one page<sup>9</sup>. There are three types of objects: *Branch Objects*, *Out-Ref-Objects (ORO)* and *Normal Objects*. A *Branch Object* decides by a random operator which object reference to follow in the tree. An *ORO* has pointers to objects in different pages which are all accessed when encountered. A *Normal Object* points to three other objects in the same page. The type for all objects has four pointers and a size of 72 bytes. In one run 195 pages are accessed and each page contains 112 objects.

The application starts with one root object from the first page. The *Branch Objects* decide the navigation in the page. When a reference to another page from the upper level (e.g. Pages 2 to 7) is encountered only the first object from

<sup>9</sup> Every page has the same structure.

the other page is dereferenced and then the application continues in page 1. At the lower levels (e.g. Pages 8 to 19) two pages are dereferenced with 1 object (the same as at the upper level) and in one page the application continues the navigation. Having two or three references to other pages gives us the possibility to test prefetching under strict time conditions. It also means that the program is quite I/O intensive and the savings in percentage terms will be high. In [15], [16] we presented a benchmark which was less I/O intensive.

Fig. 6 needs some explanation concerning the number of *Normal Objects*. The number of *Normal Objects* before an *Out-Ref-Object* is 15. The cost of processing 20 objects is equal to the cost of one page fetch in our system environment. Every object is fetched into memory with no computation or waiting time on the object.

#### 4.4. Performance Measurements

All the tests were made in a multi-user environment. Because we were unable to get exclusive access to the machines and network it is not possible to compare the absolute times of different figures. Nevertheless, we made the tests at a time when the workload was low. Savings in percent mean the percent saving of a prefetching version compared with a version without prefetching and multithreading (Demand).

In fig. 7a we present the results of our benchmark. Prefetch1 means one PrefetchThread supports the application; Prefetch2 means there is one PrefetchThread and one SupportThread and Prefetch3 has one PrefetchThread and two SupportThreads. Fig. 7a shows that with an increased number of prefetch threads the elapsed time of the applications is reduced. Recall from the benchmark structure that an ORO has three references to other pages, therefore Prefetch3 has the best performance because it achieves the optimal number of prefetch threads for page requests. Fig. 7b shows the savings of the prefetching versions in percent. Prefetch1 only provides a 5% improvement, compared with Prefetch3 which achieves a saving of 23%.

The effect of additional clients (Demand versions running on other machines) is shown fig. 8. Prefetch2 always performs better than Demand. At the level of 3 clients, Prefetch1 performs worse than Demand. This is because the cost of the thread management is higher than the benefit of prefetching. In general each client increases the network workload, the server processing and the work for the



disk manager. If the prefetch request is completed before client access the savings in percent would increase in a multi-client environment. If the prefetch request would be served at the same time at the server or even later, prefetching would decrease performance due to the additional costs of the thread management.

In fig. 9 we present the results of our distributed database test. Prefetching always generates additional workload for the server, so that a multi-server environment is more suitable for prefetching. For this test we split the database into two databases, each managed by one server. The servers both run on the same machine so as to have the same conditions for the hardware. Fig. 9 shows that all versions improve performance in the distributed environment.

The size of the buffer pool has an important impact on the performance of the prefetch technique. We compared performance for 10, 100 and 200 frames in the buffer pool. The update versions write just one object on the page, which causes the page to be marked dirty. The time for this test was stopped just before the commit of the transaction. Comparing both read versions in fig. 10, the Prefetch version can increase slightly the amount of saving with increased buffer size. The elapsed time of the Demand version increases whereas the elapsed time of the Prefetch version stays almost constant. The Prefetch version performs better with a larger number of buffer frames because this reduces locking of synchronization variables. The write versions behave very similarly. A higher number of available frames reduces the number of server flushes at transaction time, which has a direct effect on the response time.

In the next test (fig. 11) we stopped the time after the commit of the transaction. For the read versions the result are the same as in fig. 10, the Demand version increases slightly and the Prefetch version stays almost constant. For the write versions we created one version, called *Prefetch write*, which flushes all dirty pages at the end of the transaction sequentially and another version, called *Prefetch write mt flush* with has two FlushThreads to do the flushing in parallel. All write versions reduce elapsed time with a buffer size of 50 compared with 10, but their elapsed time increases after 50. Over a buffer size of 100 the multithreaded flush version outperforms the sequential flush version (at a buffer size of 200 the advantage of the multithreaded version is 1.23 seconds). This test proves that multithreading is not only useful for prefetching; flushing dirty pages to the server is also a suitable application for multithreading.

In fig. 12 we compared our presented prefetching technique (using a POT table) with the technique presented in [5] which prefetches only adjacent refer-

ences from the current object. The prefetch adjacent version achieved quite a surprisingly good result because of the structure of the benchmark. All three page fetches are from one object. Two prefetches can be made at the same time as the first page is prefetched. Therefore only the first page fetch has a sufficiently large penalty to arrive late.

The problem of prefetching is when prediction accuracy is low, prefetching can actually decrease performance. We used the same benchmark for this to test the effect of incorrect prefetch requests, with the only difference being that we navigate only to one page from the ORO. The other two pages can be used for incorrect prefetches. In fig. 13 *Two incorrect* means prefetching two pages incorrectly from an ORO; *One incorrect* means prefetching one incorrectly and *Correct* means optimal prefetching. The Inter-Reference Time (IRT) simulates overhead for client processing. The values of the IRTs are the number of loop iterations after every object access. The elapsed time of 3560 iterations is equal to the elapsed time of one object preparation for client access. 1675 iterations are equal to half of one object preparation and 850 iterations are equal to a quarter. Additional overhead for client processing could be produced by a pointer swizzling technique, application client processing or waiting time. *One incorrect* and *Correct* always perform better than Demand. After an IRT value of 850, *Two incorrect* can also improve performance.

The workload of the database client is important for the scheduling of the prefetch threads. If the prefetch thread is scheduled after encountering a PSO and the operating system time slice ends after sending the prefetch request to the server, prefetching can be even more successful under a high workload. Otherwise, if the prefetch thread is not scheduled before application access the prefetch request is unnecessary and produces processing overhead. In fig. 14 we varied the workload on the client workstation. A workload of 4 means that all four processors are fully utilized and the idle time is almost 0%. The Prefetch version performs well under the workload of 2.8 and even better above the workload of 5, i.e. where there is queueing for CPU resources. At the workload level around 4, i.e. just at the point where all processors are busy, the performance of the prefetch threads suffers as a result of operating system scheduling and therefore prefetch requests are arriving late or after the object fault.

Multithreading on the database client side can be used not only for I/O but also for very expensive CPU functions. On analyzing the client code we found out that there is an expensive function to calculate the free slot space in the

EXODUS client software. This function is called on every object access and then calculates the free slot space of the whole page. We created another thread for this function and the result of the performance test is presented in fig. 15 and fig. 16. In the test of fig. 15 we varied the number of objects accessed in a leaf page<sup>10</sup>. With an increasing number of objects in access the Audit application speeds up. All the pages that have to be checked by the Audit Thread are put into a queue and then the application thread continues with processing. If the audit page is already in the queue it is not inserted again which reduces the amount of client processing<sup>11</sup>. Fig. 16 shows that the Audit Thread has actually a higher number of buffer misses because we used the same POD but the total amount of overhead is less.

Increasing the number of SupportThreads also increases the total synchronization time of the application. We analyzed an application with one SupportThread (table 3) and an application with two SupportThreads (table 4) using the Solaris Thread Analyzer [20]. The first four rows show the I/O per second. Most of the I/O is done by the prefetch threads. In the 1 SupportThread application, the I/O work is divided over 2 threads and in the 2 SupportThreads application over 3 threads. The PrefetchThread waits on a condition variable when there is no work to do. The highest increase in synchronization time was caused by the mutex wait time. Two SupportThreads have a total mutex wait time of 7.05 seconds whereas the 1 SupportThread version only requires 5.02 seconds. When the SupportThreads are idle they are waiting on a semaphore.

For the last test we created a benchmark in which every ORO has 7 references to other pages. This benchmark was designed to test the scalability of SupportThreads. We varied the number of SupportThreads from 0 to 7. If the thread overhead is low and all the SupportThreads are scheduled in time by the operating system, the best result could be achieved with 7 SupportThread. On the other hand, if the synchronization cost of the threads is high and the SupportThreads are scheduled late then best performance is achieved by about 2 SupportThreads<sup>13</sup>. We created 4 prefetching versions which start the prefetch

<sup>10</sup> Recall the structure of the benchmark in which we accessed only one object in a page and then followed the navigation through other pages. We define such a page as a leaf page.

<sup>11</sup> To ensure the integrity of the database pages the Audit Thread must be finished before the commit of the transaction.

<sup>12</sup> CV means condition variable.

<sup>13</sup> All four threads (App, Prefetch and 2 Support) run on four processors.

Table 3  
Performance Characteristics of a 1 SupportThread application

	Prefetch T.	Support T.	App T.
file reads (bytes)	956,768	643,344	8,248
file reads (ops)	627	412	5
file write (bytes)	6,496	4,368	56
file write (ops)	116	78	1
CPU time	34.01	34.24	31.23
CV <sup>12</sup> wait time	163.23	0.00	0.00
mutex wait time	3.39	1.61	0.02
read wait time	0.56	0.40	0.00
semaphore wait time	0.70	184.89	0.00
total sync wait time	167.32	186.50	0.02

Table 4  
Performance Characteristics of a 2 SupportThreads application

	PT.	ST. 1	ST. 2	AT.
file reads (bytes)	643,344	486,632	470,136	8,248
file reads (ops)	385	302	297	6
file write (bytes)	4,368	3,304	3,192	56
file write (ops)	78	59	57	1
CPU time	35.75	35.99	35.94	32.67
CV wait time	170.88	0.00	0.00	0.00
mutex wait time	2.57	2.29	2.17	0.02
read wait time	0.28	0.19	0.20	0.00
semaphore wait time	0.74	182.33	180.53	0.00
total sync wait time	174.19	184.61	182.70	0.02

operation with different PODs. We used the values of 20, 30, 40 and 50 as a POD. Fig. 17 shows the result of this benchmark. All prefetching applications show the highest decrease from 0 to 1 SupportThread. The best result is achieved at the level of 2 SupportThreads because all threads are executed on the same processor without any context switches. After the level of 2 all applications increase elapsed time.

## 5. Conclusions and Future Work

We presented a prefetching technique for object-oriented databases using multithreading. At run time we observe the application client processing and if

the application moves towards a page which is not resident we will prefetch this page using the information of the POT. All information of the POT is collected off-line by an analyzer. This information is used at run time by the prefetch threads. We also use multithreading for flushing dirty pages to the server. This is especially useful at the end of a transaction when many pages have to be flushed to the server. The results of our benchmark showed that prefetching can improve performance significantly if object access is reasonably predictable. The fastest prefetching version achieved a saving of 23%. Even with a prefetch accuracy of only 33%, performance can still be improved by prefetching. Prefetching is also successful with additional database clients as long as the server does not become a total bottleneck. We showed that the organisation of distributed databases is attractive for prefetching. Increasing the buffer size improved significantly the response time in a write transaction. To speedup the flush of dirty pages to the server we also used threads. Above a buffer size of 100 frames the multithreaded version has a 15% percent advantage over the sequential flush version.

In the future we will compare our technique with other proposed prefetching techniques. We will implement the OO7 benchmark to see how our technique will perform in this environment. We also look at real application structures for prefetching. Another possibility is to make the ESM server multithreaded. The problem is the synchronization of the threads to access global data concurrently and safely. For example, if many threads want to access the buffer pool, waiting time will be increased. If the client demands multiple pages without strict time restrictions, our system will test the unit of I/O with a set of pages.

## Acknowledgments

We would like to thank Edward Heal, Paul Coe, Thomas Zurek, Isabel Rojas-Mujica, Steve Cusack, Tim Hopkins, Marcio Fernandez from the University of Edinburgh and Quintin Cutts from the University of Glasgow for all their comments and suggestions. Special thanks to Peter Thanisch and Stephanie Main for all their support.

## References

- [1] J.-H. Ahn and H.-J. Kim. SEOF: An Adaptable Object Prefetch Policy For Object-Oriented Database Systems. In *Proc. of the 13th Int. Conf. on Data Engineering*, Birmingham, UK,

April 1997.

- [2] P. Cao, E.W. Felten, A.R. Karlin, and L. Li. Implementation and Performance of Integrated Application-Controlled File Caching, Prefetching, and Disk Scheduling. *ACM Transactions on Computer Systems*, 14(4), November 1996.
- [3] M.J. Carey, D.J. DeWitt, G. Graefe, D.M. Haight, J.E. Richardson, D.T. Schuh, E.J. Shekita, and S.L. Vandenberg. The EXODUS Extensible DBMS Project: An Overview. In S.B. Zdonik and D. Maier, editors, *Readings in Object-Oriented Database Systems*, pages 474–499. Morgan Kaufmann, 1990.
- [4] M.J. Carey, D.J. DeWitt, and J.F. Naughton. The OO7 Benchmark. In *Proc. of the ACM SIGMOD Int. Conf. on Management of Data*, pages 12–21, Washington, USA, May 1993.
- [5] E.E. Chang and R.H. Katz. Exploiting Inheritance and Structure Semantics for Effective Clustering and Buffering in an Object-Oriented DBMS. In *Proc. of the ACM SIGMOD Conference on the Management of Data*, pages 348–357, Portland, Oregon, June 1989.
- [6] J.R. Cheng and A.R. Hurson. On the Performance Issues of Object-Based Buffering. In *Proc. First Int. Conf. on Parallel and Distributed Information System*, pages 30–37, Miami Beach, Florida, December 1991.
- [7] M.S. Day. *Client Cache Management in a Distributed Object Database*. PhD thesis, Massachusetts Institute of Technology, Laboratory for Computer Science, 1995.
- [8] D.J. DeWitt, J.F. Naughton, J.C. Shafer, and S. Venkataraman. Parallelizing OODBMS traversals: a performance evaluation. *VLDB Journal*, 5(1):3–18, January 1996.
- [9] C.S. Freedman and D.J. DeWitt. The SPIFFI Scalable Video-on-Demand System. In *Proc. of the ACM SIGMOD/PODS95 Joint Conf. on Management of Data*, pages 352–363, San Jose, CA, May 1995.
- [10] C.A. Gerlhof and A. Kemper. A Multi-Threaded Architecture for Prefetching in Object Bases. In *Proc. of the Int. Conf. on Extending Database Technology*, pages 351–364, Cambridge, UK, March 1994.
- [11] C.A. Gerlhof and A. Kemper. Prefetch Support Relations in Object Bases. In *Proc. of the Sixth Int. Workshop on Persistent Object Systems*, pages 115–126, Tarascon, Provence, France, September 1994.
- [12] J. Gray. What Happens When Processors Are Infinitely Fast and Storage Is Free? In *Fifth Workshop on I/O in Parallel and Distributed Systems (IOPADS'97)*, San Jose, CA, November 1997.
- [13] J. Karpovich, A. Grimshaw, and J. French. Extensible File Systems (ELFS) An Object-Oriented Approach to High Performance File I/O. In *Ninth Annual Conference on Object-Oriented Programming Systems, Languages, and Applications*, pages 191–204, Portland, Oregon, October 1994.
- [14] T. Keller, G. Graefe, and D. Maier. Efficient Assembly of Complex Objects. In *Proc. of the ACM SIGMOD Int. Conf. on Management of Data*, pages 148–157, Denver, USA, May 1991.
- [15] N. Knafla. A Prefetching Technique for Object-Oriented Databases. Technical Report ECS-CSG-28-97, Department of Computer Science, University of Edinburgh, January 1997.
- [16] N. Knafla. A Prefetching Technique for Object-Oriented Databases. In C. Small,

- P. Douglas, R. Johnson, P. King, and N. Martin, editors, *Advances in Databases, 15th British National Conference on Databases, BNCOD 15*, Lecture Notes in Computer Science, pages 154–168, London, United Kingdom, July 1997. Springer-Verlag.
- [17] M. Palmer and S.B. Zdonik. Fido: A Cache That Learns to Fetch. In *Proc. of the 17th Int. Conf. on Very Large Data Bases*, pages 255–264, Barcelona, Spain, September 1991.
- [18] A.J. Smith. Sequentiality and Prefetching in Database Systems. *ACM Transactions on Database Systems*, 3(3):223–247, September 1978.
- [19] SunSoft, 2550 Garcia Avenue, Mountain View, CA 94043, USA. *SunOS 5.3 Guide to Multithread Programming*, November 1993.
- [20] SunSoft, 2550 Garcia Avenue, Mountain View, CA 94043, USA. *SPARCworks / iMPact: Tools for Multithreaded Programming*, November 1995.

