

# On the Performance of Object Clustering Techniques

Manolis M. Tsangaris and Jeffrey F. Naughton

{mt,naughton}@cs.wisc.edu

Department of Computer Sciences  
University of Wisconsin–Madison

## Abstract

We investigate the performance of some of the best-known object clustering algorithms on four different workloads based upon the Tektronix benchmark. For all four workloads, stochastic clustering gave the best performance for a variety of performance metrics. Since stochastic clustering is computationally expensive, it is interesting that for every workload there was at least one cheaper clustering algorithm that matched or almost matched stochastic clustering. Unfortunately, for each workload, the algorithm that approximated stochastic clustering was different. Our experiments also demonstrated that even when the workload and object graph are fixed, the choice of the clustering algorithm depends upon the goals of the system. For example, if the goal is to perform well on traversals of small portions of the database starting with a cold cache, the important metric is the per-traversal expansion factor, and a well-chosen placement tree will be nearly optimal; if the goal is to achieve a high steady-state performance with a reasonably large cache, the appropriate metric is the number of pages to which the clustering algorithm maps the active portion of the database. For this metric, the PRP clustering algorithm, which only uses access probabilities achieves nearly optimal performance.

## 1 Introduction

In recent years a number of clustering algorithms for object-oriented databases have appeared in the literature. These algorithms attempt to improve the performance of object-oriented database systems by placing on the same page related sets of objects, thus attempting to avoid the performance penalty of one disk I/O per

object access. For the most part, these algorithms have each been presented in isolation, with some experimental data illustrating how these algorithms perform when compared to no clustering or “random” clustering. In this paper we investigate the relative performances of a number of these clustering algorithms on four different workloads based upon the Tektronix [And+90] benchmark. Our results apply to object bases with similar object structure, properties, and usage as in the Tektronix Benchmark.

The algorithms we compared were BFS, DFS, and WDFS [Sta84], Placement Trees [BD90], Cactis [DK90], PRP [YW73] and stochastic clustering [TN91a]. Of these algorithms, BFS and DFS depend only on the structure of the object graph, while the other algorithms depend in addition on a information gleaned from a training trace representative of some workload. In more detail, these algorithms are “trained” by letting them gather statistics from a trace representative of the workload; they then use these statistics and the structure of the object graph to decide upon a good clustering. To evaluate the quality of the resulting clustering, one runs another trace, different from the training trace yet still representative of the given workload, and gathers statistics about page fault rates and numbers of pages touched.

We found that for all four workloads tested, stochastic clustering gave uniformly the best results by a variety of performance metrics. Stochastic clustering works by postulating that the workload is generated by some stochastic process, then gathering statistics from the training trace to estimate the parameters of this hypothetical stochastic process, and finally mapping objects to pages so as to minimize the probability that a pair of consecutive object accesses in the reference stream crosses a page boundary. The results of these experiments were an important confirmation of the utility of the ideas behind stochastic clustering, since before performing these experiments, it was not obvious that stochastic clustering would perform this well. In particular, a number of the assumptions made by stochastic clustering are only approximately true — references in the workloads in the Tektronix benchmark are not generated by stochastic processes, and it was not im-

---

This work was supported by NSF grant number IRI-9157357

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

1992 ACM SIGMOD - 6/92/CA, USA

© 1992 ACM 0-89791-522-4/92/0005/0144...\$1.50

mediately obvious that minimizing the probability that consecutive object accesses cross page boundaries maximizes performance.

Stochastic clustering, while highly effective in these experiments, is prohibitively computationally expensive to be applied directly in many situations. In view of this fact, it is important to find lower-cost algorithms that approximate the performance of stochastic clustering. Our results were encouraging in that for each workload tested, there was at least one computationally less expensive algorithm that approximated the performance of stochastic clustering. However, unfortunately the algorithm that approximated stochastic clustering was different for each workload. This suggests that a practical clustering strategy may be a set of clustering strategies, each appropriate for a different class of workload, rather than a single monolithic strategy.

Another fact that became clear in our experiments is that even if you fix the object base and the workload, which clustering algorithm is best depends in an important way on the performance goals of the system. For example, if the goal of the system is to perform well on traversals of a small portion of the database starting with a cold cache, the important metric is the ratio of the number of pages a traversal touches to the smallest number of pages in which the objects touched by the traversal could be stored. On the other hand, if the goal of the system is to perform well in steady state with a fairly large cache, the important metric is to how many pages the clustering algorithm maps the active portion of the database. An algorithm that performs well by one metric will not necessarily perform well by the second. A particularly interesting result is that for the large cache, steady state case, the PRP (Probability Ranking Partitioning) algorithm is nearly optimal. This is surprising since the PRP algorithm makes no use of the object graph at all, clustering solely on the basis of statistics gathered from the training trace.

A final result of this study is that like high performance race horses, high performance clustering algorithms can be temperamental. That is, the “bad” clustering algorithms are relatively insensitive to differences between the training trace and the testing trace, whereas the “best” clustering algorithms show dramatic drops in performance when the testing workload contains elements of a workload that was not included in the training trace. This suggests that the highest performing algorithms may not be desirable if the reference patterns in the system vary markedly over small intervals of time.

The remainder of this paper is structured as follows: in Section 2 we describe our model of an OODBMS and the clustering problem that arises in this model. Section 3 describes our simulation environment, which we used to run the experiments. In Section 4 we present the results of our experiments, and we conclude in Section 5.

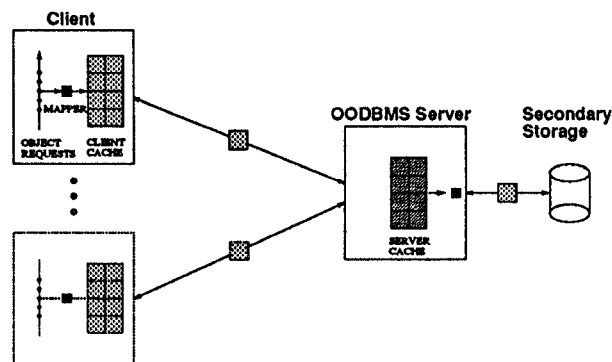


Figure 1: A simplified OODBMS Architecture

## 2 Clustering in OODBMS

A general Object Oriented Data Base Management System architecture is shown in Figure 1. The system consists of a) a data base server serving one or more clients, b) the system’s secondary memory only accessible to the server, and c) the clients, programs written in some object oriented programming language the OODBMS supports. The server supports some object abstraction for its clients, and every object in the system is uniquely identified by its Object Identifier (OID), a permanently assigned number. The object state representation is stored in the system secondary memory consisting of fixed uniform size pages. Each object is assigned to only one page through the *Clustering Mapping*. Both the server and the clients have a mechanism to find the page of an object given its OID.

Conceptually, objects can be considered vertices of a directed and possibly cyclic graph called the Object Graph. The directed edges of the object graph are just the object to object pointers. Typically, the client programs access objects *sequentially* during their execution (i.e. one at the time), by dereferencing object pointers and thus “walking on the object graph”. The program requires access to the object representation, and therefore, the client process is suspended if the object is not available. The process is resumed when the missing object is obtained from the server. In the *Paged OODBMS Architectures* the server will not just return a single object, but will return all objects mapped to the same page as the requested object. After the client receives the page containing the requested object, the client resumes execution and can access all of the objects brought in with that page.

Because of the high overhead associated with servicing requests, most OODBMS’s add caches on the client and on the server. The client cache reduces server requests, and the server cache reduces disk accesses. It is important to note that a page of the client cache contains the same objects as the corresponding page on the server. Redistributing objects to pages is usually avoided, because it can cause performance problems when writing back modified objects, and because it may

force the client to acquire many fine granularity object locks instead of fewer coarse page locks.

The choice of placing objects to pages affects the performance of the system in terms of system load, server overhead, concurrency control, and recovery. For example, 100 different objects can be placed into as many as 100 different pages, or as little as 5 pages of 4k bytes each (assuming a 200 byte average object size). The latter clustering mapping will require 1/20 as many client server interactions, 1/20 as much client cache memory, 1/20 as many pages that might need logging during updates, 1/20 as many page locks that have to be obtained, and smaller probability of conflict between different transactions. In this paper we will concentrate on the effect of clustering on memory utilization and system load.

## 2.1 Clustering Performance Measures

As the previous discussion motivates, the “packing capability” of clustering algorithms is a simple way to measure their performance. If we view objects as records and client requests as queries of records, the *Expansion Factor (EF)* can be used as such a metric. When the client requests a set of objects  $Q$  that maps to  $N(Q)$  distinct pages,  $EF$  is defined as:

$$EF(Q) = \frac{N(Q)}{\left\lceil \frac{\|Q\|}{L} \right\rceil}$$

where the denominator is the size of ideal packing of  $\|Q\|$  objects to pages of  $L$  objects each. In the above example,  $EF$  can be as low as 1 and as high as 20, but in general  $EF$  ranges between 1 and  $L$ . If the client may issue any one of  $Q_1, Q_2, \dots, Q_r$  queries, each with probability  $P(Q_i)$ , it makes sense to define the average  $EF$ :

$$\overline{EF} = \sum_{i=1}^r P(Q_i) EF(Q_i)$$

For a given set of queries and their probabilities, the ideal clustering mapping minimizes  $\overline{EF}$ . Although it is easy to avoid bad clustering mappings of  $\overline{EF} = L$ , in general it is very hard to find a clustering mapping of  $\overline{EF} = 1$  [YSLS85].

Although this definition makes sense for records and queries,  $EF$  alone is not an adequate metric for clustering. The  $EF$  measures the distribution of objects to pages, but does not take into account the order and frequency with which each object is needed. There are a large number of possible clustering mappings that have the same  $EF$ . However, not all of them achieve the same performance on small cache sizes. For example, suppose we map 6 objects  $\{a, b, c, d, e, f\}$  to 3 pages of size 2:

$$c_1 = \{[a, b], [c, d], [e, f]\} \quad c_2 = \{[a, f], [b, d], [c, e]\}$$

Both mappings  $c_1, c_2$  achieve the same  $EF = 1$  for all queries involving all 6 objects. Both perform well

with caches of size 3 pages or more. However on a cache of 2 pages they may perform very differently. A query using objects in the following sequence:  $t = abc \ abc \ abc \ \dots \ def \ def \ def \ \dots$  will thrash under  $c_2$ , but it will work fine under  $c_1$  only producing 3 page faults. It is easy to see that  $c_1$  outperforms  $c_2$  because it maps the frequently needed objects  $\{a, b, c\}$  and  $\{d, e, f\}$  to 2 pages whereas  $c_2$  maps them to 3 pages.

In general, suppose we are given a clustering mapping such that a set of objects  $Q$  maps to  $N = N(Q)$  pages. If the client cache is at least  $N$  pages big, the client will only experience an initial startup delay proportional to  $N$ , and after that, its computation will proceed at full speed. However, if the cache cannot hold  $N$  pages “thrashing” (a series of page faults) will occur, especially if some of  $Q$  objects are requested repeatedly.

It is well known that the performance of caches depends on the “locality” of page requests. Clustering maps many objects to each page and thus it increases the page locality. On small caches mappings that achieve better locality will perform better. The average working set size can be used to measure locality [Den68], the lower the working set size the higher the locality. For a sequence of requests  $(X_n) = x_1, x_2, \dots, x_n$ , the working set at time  $t$  is defined as the cardinality of the set of the last  $w$  requests.

By representing access patterns as stochastic processes, clustering can be formulated as an optimization problem [TN91a], and optimal clustering corresponds to the clustering mapping that minimizes the expected working set size. The complexity of “optimal clustering” in general (which has been shown to be NP-complete), makes it very hard to find the optimal clustering mapping. However, this formulation gives a new view to the problem, helps to come up with some more reasonable (close to “optimal”) clustering algorithms, and also gives a practical limit to how much clustering can help the performance.

The new formulation of clustering uses Markov processes to model “access patterns.” A Markov process is a sequence of correlated random variables, which in our case models object requests; the probability to request some object  $y$  at time  $t$  strictly depends on the last object requested:  $Pr\{X_{t+1} = y | X_t = x\} = P(x, y)$ . As a result, for a population of  $n$  accessible objects, this access pattern model requires at most  $n^2$  parameters the values of the  $P(x, y)$  matrix. Those parameters can be estimated from *sample sequences*, and in practice  $P$  is sparse, because of various access constraints.

## 3 Simulation Environment

In this study we have used cache simulations as a way to evaluate the performance of clustering algorithms. The evaluation of a clustering algorithm has three steps, the training, the clustering itself, and the testing. Training involves deriving a suitable access model for the algorithm (that is, a description of the access patterns) using

the object graph and/or sample traces. Next, the clustering algorithms use their access model and the object graph as input and generate a mapping of objects to fixed uniform size pages as output. Finally, the generated clustering mapping is evaluated by running test traces over the object base using a cache simulator.

### 3.1 The Workload

Our experiments were selected from the CLUB-0 clustering benchmark for the  $O_2$  system [HBD91]. CLUB-0 uses synthetic workloads based on a subset of the Tektronix Hypermodel Benchmark [And+90]. The object graph in CLUB-0 is a DAG, derived from a balanced 5-way tree with some additional edges. Figure 2 illustrates the first levels of an actual tree we have used in our simulation. Each node of the tree has edges to 10 other nodes. Five of these nodes are reachable by “child” edges, hence are called children of the parent; the other 5 nodes are reachable by “part” edges and are called subparts of the parent. The subparts of a given node are chosen at random from among all the nodes at the same level as the children of the node. The depth of the tree used in CLUB-0 is 6 resulting in 3906 objects. A query is a collection of traversals, each traversal being a sequence of object accesses representing some hypothetical operation on the graph. A traversal starts from a node high in the tree, and at each step uses a fixed rule to select the next object to visit.

Structural operations on the graph are modeled by traversals denoted as  $AXB$ , which start from a randomly selected node at level  $B$  of the graph ( $B = 0$  denotes the root of the tree).  $AnB$  denotes a Depth First Traversal (DFT) performed by exclusively following either the children hierarchy ( $1nB$ ) or the parts hierarchy ( $MnB$ ). We have added a new traversal  $AsB$  to model searches on the object graph.  $AsB$  locates some leaf level object by following a path on the graph, starting at level  $B$  and exclusively following the children hierarchy ( $1sB$ ) or the parts hierarchy ( $MsB$ ). The decision which of the parts (or children) to follow is taken randomly and independently on each level. We have also introduced *srnd*, for “skewed random”, a modified *rnd* of the Tektronix Benchmark, that visits objects randomly with a probability that follows normal distribution. The variance in the skewed workload was set to 1/10 of the object base size, and the objects were selected based on their OID but using a random permutation first. As a result, hot objects are spread “uniformly over the object base”, so there is no relationship between heat and the position in the object graph. A query is just a sequence of traversals each one starting from a randomly selected node at the starting level  $B$ . A query trace is a concatenation of its traversal traces, and contains references to objects necessary to run the query. The traversal trace is obtained by trapping all object accesses occurring during the execution of that traversal, as if the code was running on an object oriented run time system. As a guide, we used the code produced

by the non swizzling E compiler [RC89] with all optimizations having to do with persistent objects turned off, so that object references appear every time the original traversal code requires access to the object state. Finally, for our purposes a workload is a trace obtained by mixing, concatenating, or interleaving query traces.

### 3.2 Clustering Algorithms

We have implemented and tested several algorithms based on heuristics and ideas discussed in the literature. The goal of all algorithms examined is to partition the object graph (OG) by assigning objects to uniform size pages. The object graph is formed considering objects as nodes and any reference from an object to some other object as a directed edge connecting them. Most clustering algorithms use as input a graph representation of the access patterns (called clustering graph or CG), assumed to be characteristic of the client behavior. The representation is usually derived from the object graph and/or from sample traces (the training traces).

In general, three types of CGs are used:

- The OG, i.e. the object graph itself; a rudimentary representation of access patterns that does not take advantage of any other knowledge that may be available. No statistical information from the training traces is captured in OG.
- The SG (Statistical object Graph); the object graph annotated with edge and node weights. The node weight (edge weight) is equal to the frequency the object (edge) appears in the training trace.
- The SMC (Simple Markov Chain) is the directed graph form of a first order Markov process that could have produced the object trace. For each accessible object in the system the graph contains a node. Any positive probability that one object can be accessed after some other object, is represented as a directed edge. The node (edge) weights are the the estimated stationary (transition) probabilities of the chain from the sample trace.

The plain object graph does not convey much access information. The last two graphs express the behavior of the client as it is manifested in the training trace, using a much more compact representation than the trace itself. The SG limits its information in the usage of objects and references, whereas SMC records arbitrary but only consecutive transitions.

The majority of the algorithms assign objects to clusters sequentially, by performing a form of graph traversal on CG and assigning objects to clusters as they go. On each step the object is put in to the current cluster and if it fills up, a new empty cluster is created. Table 1 summarizes the algorithms used. The notation in that table is to generate the name of an algorithm from the type of clustering graph it uses (OG, SG, or SMC) and the style of the graph traversal or operation they perform.

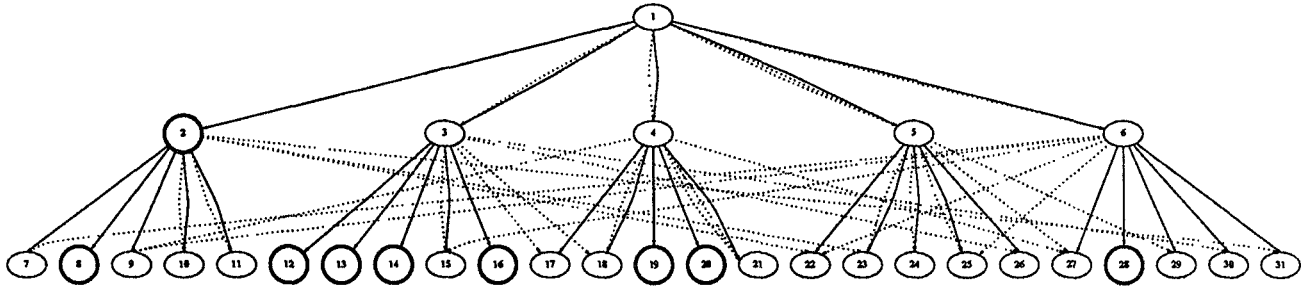


Figure 2: The Tektronix Benchmark Object Graph

The SMC.PRP SMC.KL algorithms were proposed in [TN91a]. The PRP (initially proposed for record clustering in [YW73]) method just uses the node weights of the SMC graph (i.e. the absolute probabilities), by sorting objects with respect to their probability and then assigning them to pages in that order. This scheme is also known as Probability Ranking Partitioning. The SMC.PRP algorithm has  $O(n \log n)$  cost. The SMC.KL algorithm, uses the standard Kernighan–Lin [KL70] graph partitioning algorithm to find a near to optimal clustering of the SMC graph. SMC.KL is the algorithm we have referred to in the introduction as “Stochastic Clustering.” SMC.KL partitions the object graph so that the expected working set for window size 2 is minimized. SMC.KL is a heuristic partitioning algorithm that achieves only pairwise optimality, i.e., there will be no two nodes belonging to two different partitions that can be exchanged and result in a lower total cost partitioning. SMC.KL does not cluster sequentially, since it applies repartitioning until no cost improvement is possible. The complexity of SMC.KL is dominated by the complexity of graph partitioning and it is on the average  $O(n^{2.4})$ [PS82].

SMC.WISC, a new algorithm we propose, appears in many cases to approximate well SMC.KL in terms of performance. WISC is a traversal based (greedy) low cost graph partitioning that does non-backtracking clustering. Objects are visited with the order of their absolute probabilities, hotter objects first. Each unclustered visited object is selected to start a partition. While there is room in the current partition, all objects *accessible* in terms of the SMC graph from the current contents of the partition are considered. The object that maximizes the overall probability of using that partition, is selected and the process is repeated until the partition fills up. At this point, the next unclustered object from the sorted list is considered, and the whole process is repeated, until all objects have been clustered. SMC.WISC has cost  $O(n \log(n))$ .

The OG.DFS algorithm traverses the object graph in a DFS manner. It minimizes the number of different pages touched during a pure DFT that uses all possible object graph edges. OG.BFS traverses the graph in a BFS manner, grouping siblings together as much as possible. Both algorithms have linear cost  $O(n)$ .

SG.WDFS is much like the OG.DFS except that dur-

Algorithm	Complexity	Proposed
OG.DFS	$O(n)$	[Sta84]
OG.BFS	$O(n)$	[Sta84]
OG.PT	$O(n)$	[BD90]
SG.WDFS	$O(n)$	[Sta84]
SG.CACTIS	$O(n \log n)$	[DK90]
SMC.PRP	$O(n \log n)$	[TN91a]
SMC.WISC	$O(n \log n)$	here
SMC.KL	$O(n^{2.4})$	[TN91a]

Table 1: Tested clustering algorithms

ing the DFT, siblings are selected depending on how hot their edge is. This algorithm attempts to minimize the number of clusters “probable” DFTs will touch. Edge probability information can be supplied by the compiler based on static usage information (like in Semantic Clustering [SS90]), or as user hints (like those originally planned for E [RC89]). In our case, SG.WDFS gets its hints from the sample traces. SG.WDFS has linear cost  $O(n)$ . OG.BFS, OG.DFS and SG.WDFS have been proposed first in [Sta84] for clustering Smalltalk objects.

The SG.CACTIS is based upon the clustering algorithm proposed in [HK89], a heuristic algorithm that performs PRP scan of the objects, clustering the closure of each group as it is formed.

Finally, we have implemented OG.PT, an algorithm based on placement trees. A placement tree is a pattern that matches a subset of objects reachable from a given node, the root of the placement tree. The matched objects are always connected with object pointers and form a tree. OG.PT traverses the object graph in some manner (e.g. BFS or DFS), and matches a given placement tree against the object graph starting at the visited object. All matched un-clustered objects are inserted to a logical cluster, which is subsequently assigned to physical pages. Although OG.PT is very intuitive, it is not an automatic method. On  $O_2$  the data base administrator selects the placement trees [Deu+91] based on his system experience. In our implementation, OG.PT refers to the performance of the best placement tree we were able to find for the workload in question.

### 3.3 Cache Simulation and Performance Metrics

Each clustering mapping was tested using a client simulator as in [TN91a], and [HBD91]. The input to the simulator is the testing trace and a clustering mapping. Each object reference is converted to the corresponding page reference using the mapping. Then, that page is retrieved from a variable size LRU cache, and the number of page hits is updated for each cache size. Periodically or at the end of the simulation the average cache hit rate is reported.

The primary performance metrics used in this study, are the client cache hit ratio  $HR$  and the expansion factor  $EF$  (defined in Section 2).  $EF$  is more appropriate when cache sizes are large, where the particular arrangement of objects to pages does not matter.  $EF$  is less meaningful when client caches are small compared to the portion of the object base being accessed.  $HR$  illustrates better the performance of clustering mappings, when traversals are longer and fill up a small cache.

## 4 Results

In this section we present the results of our experiments with the performance of the clustering algorithms on a variety of workloads. In addition to exploring the performance of the algorithms on “pure” workloads, were testing and training workloads are statistically similar, we will present results from “noisy” and “mixed” workloads. The experiments presented here are for a uniform object sizes of 200 bytes (typical, as reported in [Bai89]), and pages of 4k bytes each.

### 4.1 Performance of single traversals

This experiment tests the algorithms with respect only to their  $EF$ , i.e. their capability to group all objects requested during a single traversal as dense as possible. We measured  $EF$  by running the client on a large cold cache, and counting the number of page faults at the end of the traversal. Since the client cache is initially empty this number gives exactly  $N(Q)$ , the number of pages the traversal is mapped to (also called “traversal pages”).  $EF$  is also an indication of the “loading time”, or how quickly the required objects for a traversal are brought in. To estimate  $\overline{EF}$ , the experiment is repeated a number of times, each time selecting a different traversal.

Figure 3 shows the  $\overline{EF}$  achieved by clustering algorithms on a variety of workloads. Note that one of the curves in the graph has three labels: OG.BFS, OG.DFS, and SG.WDFS. The meaning of this notation is that the curves for those three graphs were virtually indistinguishable within the margin of experimental error. We will follow this convention of condensing indistinguishable curves and their labels to a single curve with the union of the original labels in the graphs throughout this section. Workloads are ordered with respect

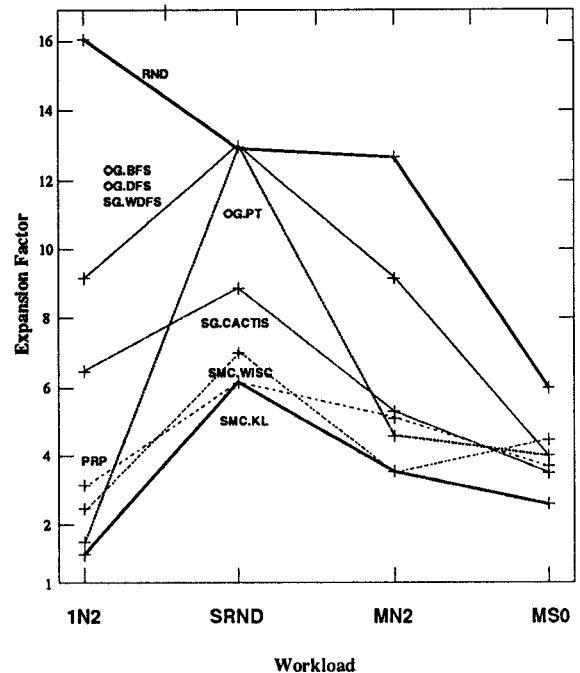


Figure 3: Performance on pure workloads

to the amount of active objects, i.e. the percentage of the object graph they use. 1n2 accesses the most objects and mn2/ms0 the least.<sup>1</sup> The expansion factor of all algorithms fluctuates as the workload changes, and random clustering gets better as fewer the objects are accessed. Statistical algorithms (SG.CACTIS, SG.WDFS, SMC.KL, SMC.WISC) are more adaptive, since they take advantage of the access pattern training. Algorithms based only on the object graph are less adaptive, since they do not train on access patterns. Most algorithms do not have a constant rank for all workloads, except of SMC.KL (best) and RND (worst). Note that the skewed random query (srnd) makes most algorithms as bad as random clustering, since srnd does not follow the object graph at all.

Notice that PRP performs well, although not optimally, on all of these workloads. The good performance of PRP on the object-graph traversal workloads (mn2, 1n2, and ms0) is surprising, since PRP never looks at the object graph and clusters based solely on the zero-order statistics from the training trace. The performance can be explained as follows. Since the root objects for the traversals in the training trace are randomly selected, due to randomness these objects will be selected with slightly differing frequency. Since every traversal out of the same root object is identical, every object (except for the leaves of the traversal) is selected with a frequency identical to that for the root of the traversal. Ignoring duplicates between traversals (which are relatively few in these workloads), this in turn means that the frequencies for objects belonging

<sup>1</sup>We will postpone the explanation to Section 4.2

to the same traversal are identical and slightly different from the frequencies for objects belonging to different traversals. Finally, this means that since PRP groups objects according to decreasing frequency of access, it will tend to cluster an object with other objects that belong to the same traversal. In effect, PRP is using the different frequencies of objects belonging to different traversals to “learn” the structure of the object graph.

Placement trees are extremely good in the 1n2 queries, where traversals are disjoint and well known in advance. Each traversal accesses its own subtree rooted at level two, and a placement tree can easily pack objects exactly that way. As a result OG.PT gets an  $EF$  close to 1. In non-graph queries like `srnd`, placement trees are not applicable at all, since by definition placement trees attempt to group together only objects connected by edges of the object graph. Since the traversals of `mn2` (`ms0`) are not disjoint (there can be multiple paths through part edges to the same node), we could not come up with placement trees that group each traversal as well as in the 1n2 case. Finally, the simple structural algorithms (OG.DFS,OG.BFS) cannot compete with the more sophisticated placement trees.

The stochastic clustering algorithm is definitely a winner in terms of  $EF$ . It performs as good as the manually chosen placement trees in its ideal case (1n2), it adapts to arbitrary random queries (`srnd`) as well as to the queries that traverse the object graph in a variety of ways.

## 4.2 Hot Cache Performance

It is intuitively known that large caches are more forgiving to less efficient clustering mappings than are small caches, an effect we want to illustrate with this experiment. If the object graph is used repeatedly and the client cache is “large”, then  $EF$  (or even  $\overline{EF}$ ) is not the right metric of clustering performance. In the case of cold caches the number of page faults is proportional to  $EF$ . If the cache is not empty, then in addition to the intra-traversal locality, the hit ratio is influenced by the amount of pages shared between traversals (the inter-traversal locality).

The cache gets “hotter” as more and more traversals are executed on the client. The first traversals achieve a lower hit ratio, and the cache is in effect loaded with pages that were used by old traversals. Subsequent traversals find more and more of their needed pages in the cache, and their hit ratio improves. A simple analysis of random clustering suggests that very quickly, the whole object base will be touched, since there is a high probability that each object referenced by a traversal will be found on a different page. In the `mn2` case each traversal makes 311 object requests and accesses up to 156 distinct objects. Under random clustering, about 1040 randomly selected objects are required to touch all 196 pages of the object graph. To touch 1040 objects, 9 different traversals are required approximately,<sup>2</sup>

<sup>2</sup>This number is only approximate because of the possible du-

thereby producing 2800 `mn2` references. Indeed, as the results of Figure 4 show, random clustering arrives at the steady state in approximately 3000 references. Random clustering might seem an ideal way to increase the inter-traversal locality, and this is certainly true if the whole object graph is accessed with the same probability and ample cache memory is available.

However, an interesting property of the Tektronix Benchmark object graph restricts the number of objects accessed during `mn2` traversals. The `mn2` query uses exclusively the parts relationship, indicated by the dashed edges of Figure 2. In our case on each level  $h$  of the tree (level 0 being the root) there are  $A(h) = 5^h$  nodes that point to  $A(h+1) = 5^{h+1}$  children through the parts edges. Since each node has only 5 parts, the expected number of orphans (i.e. children that are not connected by the parts edges) at level  $h$  is:

$$A(h) \left( 1 - \frac{1}{A(h)} \right)^{A(h)}$$

Those nodes cannot be reached by any node higher in the tree and are represented as round objects in Figure 2. If the previous formula is applied recursively from the starting level up to the leaves, it can be derived that the expected number of total nodes accessible for the `mn2` query that starts from the second level of the 6-level tree, is just 33% of the total (for the derivation please refer to [TN92]).

The graph of Figure 4 shows the cumulative number of pages accessed during `mn2` traversals as a function of the number of references, using a large initially empty cache, that can hold the whole object base. The graph begins after the first traversal has been performed, so the curves start from a position proportional to  $EF$ . Initially all clustering mappings rapidly bring in a number of pages, and after some point, they have a very slow page fault rate thereby reaching the steady state. The curve labelled as OG.KL corresponds to optimal clustering of the object graph using the Kernighan–Lin partitioning algorithm and giving all edges equal weights. In most of the cases, OG.KL performed no different than random clustering, and therefore straightforward partitioning even of a relatively sparse object graph may not be a good clustering heuristic. The OG.PT-0 curve corresponds to the second to the best placement tree (OG.PT) we found for this workload, and it performs much worse.

The number of pages touched at the steady state shows an important property of clustering algorithms; their capability to map the active portion of the database to the minimum possible number of pages. The *long term expansion factor* or  $EF_\infty$  is an indicator of the steady state performance.  $EF_\infty$  is the ratio of pages accessed in the steady state ( $N_\infty$ ) to the number of pages that would be required ideally to pack all active objects ( $n_\infty$ ).

plicates in the `mn2` traversals and the randomness in selecting distinct traversals from all 25 possible ones.

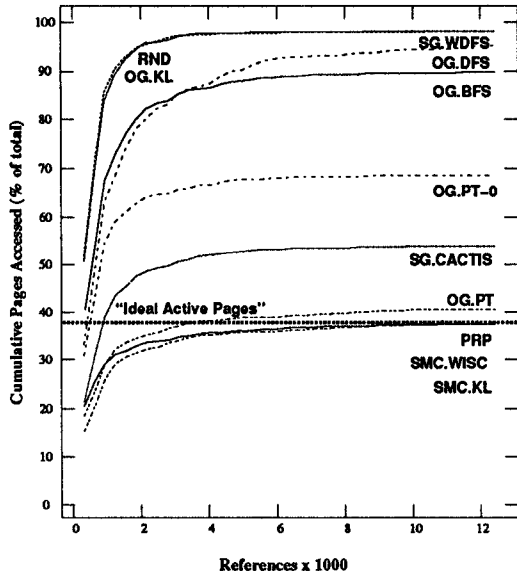


Figure 4: mn2 page faults going to steady state

In the object graph we tested, there were actually 1498 accessible objects through the parts edges, corresponding to 38% of the object graph<sup>3</sup> represented by the dashed line in the graph of Figure 4.  $EF_{\infty}$  is definitely related to  $EF$ ;  $EF$  measures the average packing capability for each one of the possible traversals that influence  $EF_{\infty}$ . In our case, since there are 25 possible equiprobable traversals,  $EF_{\infty}$  cannot be worse than  $25EF$ . However, large  $EF$  does not necessarily mean large  $EF_{\infty}$ . Imagine a (possibly unrealistic) clustering mapping, that arbitrarily packs all accessible objects to a small number of pages. This mapping may have a bad  $EF$  since active objects accessed at different traversals will be mixed, but its  $EF_{\infty}$  will be the optimal since no inaccessible objects will be clustered with accessible objects.

If only  $EF_{\infty}$  is our concern (i.e. when the system operates near the steady state), then SMC.PRP suffices. PRP uses the absolute access probabilities of objects, and packs objects according to it. As a result PRP achieves the minimum  $EF_{\infty}$ . Note that this is true even though PRP makes no use whatsoever of the object graph. The performance of placement trees in steady state can now be explained, since they manually assign all objects that can possibly be accessed during a traversal to a single cluster, assuming that in the worst case all reachable objects are used. Although their  $EF$  was higher than SMC.KL (see Figure 3), their  $EF_{\infty}$  was not much worse than optimal, due to sharing of pages between traversals.

The graph of Figure 5 is another interpretation of the steady state performance, indicating the minimum cache size needed to guarantee a given hit ratio at the steady state (note that the graph ends at 99 % and not

<sup>3</sup>It is 38% on this particular graph and not exactly 33% due to statistical variation.

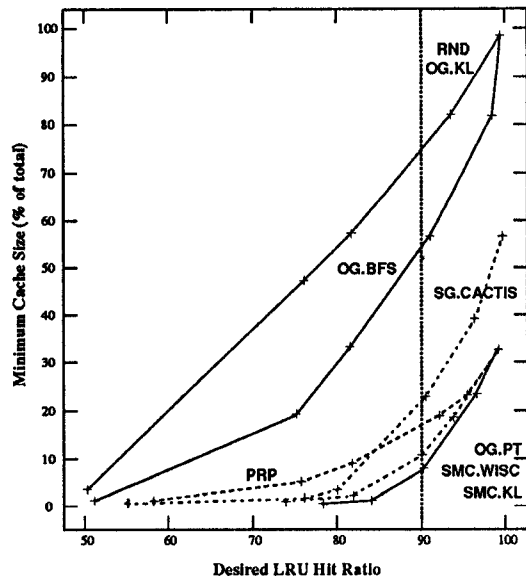


Figure 5: Min cache size to achieve a hit ratio

100 %). The graph can be used for selecting a clustering algorithm for the mn2 workload, given an amount of available memory and a desired hit ratio. If the cache size can be close to the  $EF_{\infty}$  of a clustering algorithm, then the hit ratio will be close to 100%, and the algorithm is acceptable. If memory is a constraint, then potentially more expensive algorithms with smaller  $EF_{\infty}$  should be used. For example at the steady state, random clustering we achieve 90% hit ratio by caching 80% of the object graph. The same hit ratio can be achieved with SMC.KL/WISC/OG.PT by caching less than 10% of the object base.

### 4.3 Increasing the problem size

One of the challenges in doing scaleup experiments is defining exactly what it means to scale this benchmark, since adding more levels to the object graph changes both the structural locality in the graph and also the number of objects touched per traversal. We ran tests for a variety of different ways to scale the database; we found that there were no surprises provided that the size of each traversal scales with the graph size. [TN92] contains a complete discussion of the scaleup experiments and their results, which are omitted here due to lack of space.

### 4.4 The effect of Noise

Most clustering algorithms base their performance on the knowledge of access patterns, as it is registered in their access models. Real access patterns however, may be different than the ones used for training. One way real access patterns may be different, is that some unexpected references may appear in the actual trace. To study this effect we added white noise to the testing

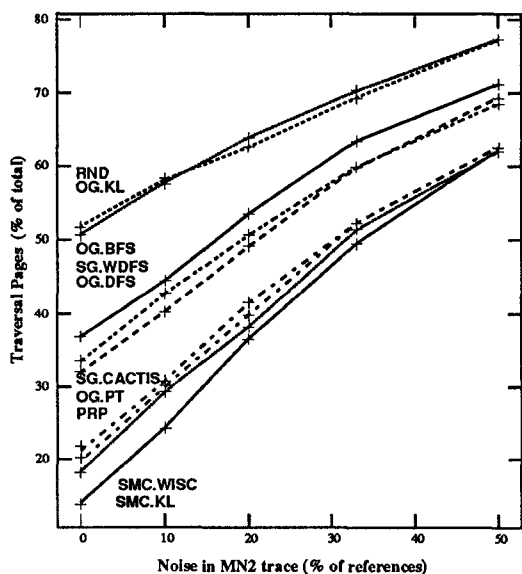


Figure 6: “Noisy” references in mn2 queries

object stream (as in [PZ91]). *White Noise* is a stream of random references chosen with uniform probability from the whole object population.

The graph of Figure 6 shows the average number of traversal pages of mn2 as a function of noise level. For small caches, practically every noisy reference is a miss, so we should expect an increase of the page faults with noise level. With 20% noise level there is one noisy reference every 5 mn2 references, resulting to 62 random references during an mn2 traversal. 62 random objects map to 53 pages (out of 200) on the average. Since SMC.KL maps an mn2 traversal to 28 pages, there will be 7 pages in common on the average and therefore, 74 pages accessed totally (or 38 %). Similarly, random maps 156 objects to 101 pages and has 27 pages in common with the noise on the average. As a result, it should require 128 pages (65 %) and indeed it does.

#### 4.5 Changing Access Patterns

So far, clustering algorithms were tested on pure traces, that is, traces containing traversals from a single query. The mixed workloads contain traversals from more than one pure workload. Mixed workloads will be used to observe the behavior of the clustering algorithms when the training trace is not entirely indicative of the testing trace.

Since efficient statistical algorithms are heavily dependent on the expected access patterns, they should be affected the most by differences between testing and training. Less efficient algorithms will be less affected, and finally random clustering should not be affected at all. This phenomenon resembles the behavior of a filter designed to accept a certain frequency and reject others, so we will call it the *tuning effect*. The interesting question is the tuning sensitivity and the next experiment

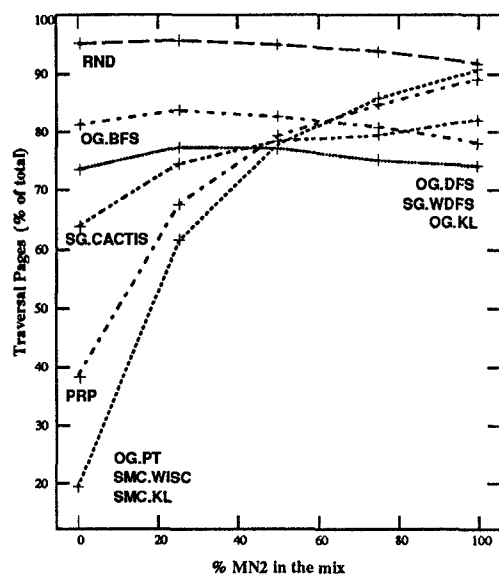


Figure 7: Changing access patterns

was done to investigate it.

In this experiment the clustering algorithms were trained with a pure 1n2 workload, but are tested with a varying mix of 1n2 and mn2 traversals. Although functionally both traversals perform the same DFT starting from the same randomly selected nodes, they follow different edges in the object graph, and generate different access patterns. The graph of Figure 7 shows the (average) number of traversal pages as a function of the mix proportion. The random algorithm remains practically unaffected, and the less efficient algorithms did not lose much either since their performance does not depend on the access patterns. The most efficient algorithms remain good initially, but they lose performance quickly and for a mix of over 50 % they are no longer the best.

From this experiment, one can conclude that the simple structural algorithms (which behave like coarse filters) are not bad for dealing with highly unpredictable access patterns. Good algorithms behave as more fine filters and great care should be taken when using them, since a significant change in the access patterns may affect them dramatically. We did additional experiments that trained on a mixed workload and tested on a varying mix; that mix did not contain “unknown access patterns” but simply different proportion of known access patterns, and the tuning effect was not as severe as here.

## 5 Conclusions

We have investigated the performance of a number of well-known clustering algorithms over a standard object-oriented benchmark, the Tektronix Hyper-model benchmark. The main points observed in our study are

that

- Stochastic clustering, while expensive, performed the best in all the tests we ran. While other, less expensive algorithms on occasion performed similarly to stochastic clustering, no single algorithm was close for all workloads, so if stochastic clustering cannot be used, care should be taken to match an appropriate inexpensive clustering algorithm with a given application. Apart from SMC.WISC a good alternative is SG.CACTIS and OG.PT although placement trees are not always easy to come up with.
- The more precise the clustering algorithm, the more sensitive it is to mismatches between training and testing access patterns.
- In contrast to mismatched access patterns, sporadic unexpected references affect the clustering performance of algorithms by approximately the same degree, and the algorithm ranking is not affected at all. As a result, such references can be safely ignored during the statistics gathering process.
- For cold-cache traversals of small portions of the object graph, the expansion factor is the important performance metric; for steady-state large cache performance the number of pages to which the clustering algorithm maps the active portion of the database is the appropriate metric. A clustering algorithm that performs well on one metric may not perform as well the other.
- Structural clustering techniques and especially placement trees, are very effective when a subset of the object graph edges is almost exclusively used to reach objects. However, when accesses are not done through the object graph (like the SRND queries), or when they use most edges of a complex graph structure (like mixed mn2/1n2 workloads), additional statistical information is necessary to produce the correct clustering mapping.

We are currently investigating further issues in object clustering, including the effect of non uniform object sizes, low-cost approximations to stochastic clustering, efficient algorithms for re-clustering sets of objects, and techniques for generating and propagating statistics throughout the database based upon type information and partial access statistics.

## References

- [And+90] T. Lougenia Anderson et al. The Tektronix HyperModel Benchmark. *EDBT*, 1990.
- [Bai89] P. Bailey. Performance evaluation in a persistent object system. In *Proceedings of the Third International Workshop on Persistent Object Systems*, Newcastle, Australia, September 1989.
- [BD90] Veronique Benzaken and Claude Delobel. Enhancing performance in a persistent object store: Clustering strategies in  $O_2$ . Technical Report 50-90, Altair, August 1990.
- [Den68] P. J. Denning. The working set model of program behavior. *Comm. ACM*, 11(5):323–333, May 1968.
- [Deu+91] O. Deux et al. The  $O_2$  system. *ACM Communications*, 34(10):34–49, 1991.
- [DK90] Pamela Drew and Roger King. The performance and utility of the CACTIS implementation algorithms. In *Proceedings of the 16-th VLDB Conference*, pages 135–147, Brisbane, Australia, 1990.
- [HBD91] Gilbert Harrus, Veronique Benzaken, and Claude Delobel. Measuring performance of clustering strategies: the CluB-0 benchmark. Technical Report 66-91, Altair, January 1991.
- [HK89] Scott E. Hudson and Roger King. Cactis: A self-adaptive, concurrent implementation of an object-oriented database management system. *ACM Transactions on Data Base Systems*, 14(3):291–321, September 1989.
- [KL70] B.W. Kernighan and S. Lin. An efficient heuristic procedure for partitioning graphs. *Bell System Technical Journal*, 49(2):291–307, February 1970.
- [PS82] Christos H. Papadimitriou and Kenneth Steiglitz. *Combinatorial Optimization: Algorithms and Complexity*. Prentice-Hall Inc, 1982.
- [PZ91] Mark Palmer and Stanley B. Zdonik. FIDO: a cache that learns to fetch. In *Proceedings of the 17-th VLDB Conference*, Barcelona Spain, 1991.
- [RC89] J. Richardson and M. Carey. Persistence in the E Language: Issues and Implementation. *Software Practice and Experience*, 19, 12 1989.
- [SS90] Karen Shannon and Richard Snodgrass. Semantic clustering. In *Proceedings of the 4-th Int'l Workshop in Persistent Object Systems*, pages 361–374, Martha's Vineyard, MA, September 1990.
- [Sta84] James W. Stamos. Static grouping of small objects to enhance performance of a paged virtual memory. *ACM Transactions on Computer Systems*, 2(2):155–180, May 1984.
- [TN91a] Manolis M. Tsangaris and Jeffrey F. Naughton. A stochastic approach for clustering in object stores. In *Proceedings of the SIGMOD International Conference on Management of Data*, pages 12–21, Denver, Colorado, May 1991.
- [TN91b] Manolis M. Tsangaris and Jeffrey F. Naughton. A stochastic approach for clustering in object stores. Technical Report 1017, University of Wisconsin, Computer Sciences Dept., April 1991.
- [TN92] Manolis M. Tsangaris and Jeffrey F. Naughton. On the performance of object clustering techniques. Technical report, University of Wisconsin, Computer Sciences Dept., March 1992.
- [YSL85] C. T. Yu, Cheing-Mei Suen, K. Lam, and M. K. Siu. Adaptive record clustering. *ACM Transactions on Data Base Systems*, 10(2):180–204, June 1985.
- [YW73] P.C. Yue and C.K. Wong. On the optimality of the probability ranking scheme in storage applications. *JACM*, 20(4):624–633, October 1973.