

# A Case for Fractured Mirrors

Ravishankar Ramamurthy David J. DeWitt Qi Su

Department of Computer Sciences  
University of Wisconsin-Madison

ravi@cs.wisc.edu, dewitt@cs.wisc.edu, qi@cs.wisc.edu

## Abstract

The Decomposition Storage Model (DSM) vertically partitions all attributes of a given relation. DSM has excellent I/O behavior when the number of attributes touched in the query is small. It also has a better cache footprint than the N-ary storage model (NSM) that is used by most database systems. However, DSM incurs a high cost in reconstructing the original tuple from the partitions. We first revisit some of the performance problems associated with DSM. We suggest a simple indexing strategy and compare different reconstruction algorithms. The paper then proposes a new mirroring scheme, termed **fractured mirrors**, using both NSM and DSM models. This scheme combines the best aspects of both models, along with the added benefit of mirroring to better serve an ad-hoc query workload. A prototype system has been built using the Shore storage manager and performance is evaluated using queries from the TPC-H workload.

## 1. Introduction

A number of the fundamental assumptions upon which the current generation of database systems are based have changed dramatically over the past decade. CPU speeds are improving rapidly (recently even faster than Moore's law would have predicted) and the amount of main memory that is affordable is also increasing. While disk capacities have also shown similar improvements, disk seek times and effective transfer rates (transfer rate/capacity) have improved at a much slower rate (almost a factor of 10 slower). In addition since it appears that disk capacities are growing faster than database sizes, even the benefits of using parallelism are likely to diminish. Hence

disk I/O will certainly constitute the primary performance bottleneck. Moreover, in modern architectures, cache performance has also been shown to be an important factor in the CPU time of query execution [11]. Hence database storage architectures that are more conscious of disk-arm optimizations and cache effects during query processing are needed.

Database systems usually store all the attributes of a relation together. But this format is not ideal for modern database architectures given that cache-misses form an important component of query execution time [1]. An alternate storage model (DSM) uses vertically partitioned tables [8]. In this representation, each attribute of a relation is stored as a separate relation along with a surrogate that identifies the original tuple that the attribute came from.

A	B	C	ID	A	ID	B	ID	C
A1	B1	C1	1	A1	1	B1	1	C1
A2	B2	C2	2	A2	2	B2	2	C2
A3	B3	C3	3	A3	3	B3	3	C3
A4	B4	C4	4	A4	4	B4	4	C4
A5	B5	C5	5	A5	5	B5	5	C5

The figure above shows a sample relation in the NSM representation on the far left and the corresponding DSM representation on the right. As described in [8], the DSM model maintains two copies of each partition, one clustered on IDs as shown above, and the second clustered on the attribute value, which serves as an index. DSM seems to have good I/O behavior when the number of attributes touched by a query is low. Consider a sample scenario in which a selection operation has low projectivity and low selectivity, i.e. only a few attributes are projected from a large percentage of the tuples. With the DSM representation only the partitions required by the query would be scanned, minimizing the number of disk I/Os performed while maximizing L1 and L2 data cache performance. With the NSM representation, since the query predicate is not very selective, an index would not be useful and the entire relation would be scanned. In addition, NSM would have poor cache performance [1].

---

*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 VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment*

PAX is a recently proposed alternative implementation of the NSM representation that employs vertical partitioning within each page [2]. For this example query, PAX would have a much better cache footprint than NSM, while having the same I/O characteristics as NSM. Hence for this query, it seems that DSM is the best choice. However, it is just as easy to come up with examples where the NSM representation is better. While DSM seems to be ideal for selections with low projectivity and low selectivity, as the projectivity increases, the cost of reconstructing the original tuple from the partitions begins to dominate the execution time. On the other hand, the NSM model is tuned for workloads that are highly selective and uses most of the attributes. Hence, neither storage format is optimal for all queries. The paper proposes a storage architecture that is a variant of the existing mirroring technique as a first step towards addressing this problem.

Mirroring [5][14] (RAID 1) is a technique for providing fault tolerance that maintains two (or more) identical copies of each disk. If one disk of the mirrored pair fails, the system can continue operating while the failed disk is replaced and then recovered from the mirror. Mirroring can be implemented in either hardware or software. In addition to providing fault tolerance, mirrors can also be used to improve performance by partitioning random seeks across the mirrored pair [5]. This can be critical since random seeks are very slow.

In this paper we propose a new form of mirroring termed fractured mirrors. With this scheme, instead of the two disks in the mirrored pair being physically identical, they are logically identical. The naïve implementation of fractured mirrors would store the NSM copy of a table on one disk of the mirrored pair and the DSM copy on the second disk<sup>1</sup>.

This scheme retains the advantages of both the NSM and DSM representations. Queries touching only a few attributes of a large number of rows will use the DSM copy. Highly selective queries or queries requiring a majority of the attributes will use the NSM copy. This idea builds on the idea of disk shadowing [5][14], which demonstrated that mirrors could be used profitably during query processing and not only for the purposes of fault tolerance. By storing the mirrors in different storage formats, we can formulate query plans that can truly maximize disk utilization while minimizing the number of L1 and L2 cache misses. However a naïve implementation of DSM can lead to surprisingly bad performance, even when very few attributes are used. The next section

<sup>1</sup> If a table is horizontally striped across multiple mirrored pairs, the rows stored on a single mirrored pair will be stored in their NSM representation on one disk and their DSM representation on the second.

details the performance limitations of the DSM representation and how simple storage schemes and scan algorithms can improve its limitations. Section 3 describes our prototype, strategies for structuring the fractured mirrors and alternative query execution plans that the use of fractured mirrors can provide. An experimental evaluation of the prototype using queries from the TPC-H benchmark suite follows. The paper concludes with a presentation of related work and strategies for handling updates.

## 2. Storing and Scanning DSM Partitions

### 2.1 The Naive Implementation

The straightforward way of implementing the vertical partitions of the DSM model is to store each vertical partition as a separate relation with two attributes, an integer that acts as an identifier and the column's value as illustrated previously. When used to store the Line-item table from the 1 GB version of TPC-H benchmark (which has 16 attributes) this approach has very poor space utilization and performance. With the NSM representation, this table occupies about 1.1 GB and a full table scan takes 74.5 seconds.<sup>2</sup> The DSM representation on the other hand occupies 2.8 GB. For DSM, the attributes were assembled one tuple at a time like a traditional scan. The following table illustrates that a naïve implementation of DSM provides a performance advantage only when a single attribute from the table is touched.

Projectivity	1	2	4	8
ScanTime(s)	68.29	138.06	366.86	759.39

Several factors contribute to the poor performance of this implementation strategy. First, the naïve DSM implementation stores each (ID, <AttrValue>) pair as standard database record on a slotted page [15]. While the slot overhead is generally not significant when a record is used to hold a tuple in the NSM representation, it can become significant when the record is used for a single (ID, <Attr Value>) pair. Furthermore, for fixed-length attributes, whose position on the page can be computed from the length of the attribute and the ID, the slotted page representation is redundant. The second significant source of wasted space is the ID itself. Storing a 32 or 64 bit identifier with each attribute can easily double the space required to store a table. It is very important to keep in mind that the real issue is not the disk space required for the slot array entry or the ID. Disk space is almost free these days. The issue is that the extra space consumes

<sup>2</sup> The NSM copy and DSM partitions were stored as files in Shore configured to have a 128 MB Buffer pool and 32 KB page size.

precious I/O disk operations when the partition is processed.

Another drawback of the naïve strategy is that it is not possible to quickly reassemble a tuple from its vertical partitions given the tuple's ID. Some form of index such as a B-Tree mapping ID to <AttrValue> is required to do this efficiently.

This leads to an alternative representation in which each vertical partition is stored as a B-Tree on ID with the leaf pages containing (ID, <Attr Value>) pairs. While this approach still wastes space storing a tuple's ID once for each of attribute values plus the cost of a slot array entry, it makes the task of reassembling a tuple given its ID straightforward. More importantly, it leads us to a refined representation that we describe below.

## 2.2 A Sparse B-Tree Based Representation

Our refined design uses a modified B-Tree design in which the overhead of the redundant IDs is eliminated for both fixed and variable length attributes and the slot array overhead is eliminated for fixed length attributes.

Our approach is based on several simple observations. First, IDs are system generated by incrementing a counter, and are never reused<sup>3</sup>. Thus, new (ID, <Attr Value>) pairs are always appended to the right-most leaf node of the B-Tree. In addition, for fixed-length attributes, given the ID of the lowest attribute value on a leaf page, there is no need to store the IDs of the remaining attributes as they can be computed given the attribute's offset from the start of the page. This avoids the need for either a slot-array or IDs. Furthermore, if storage extents are guaranteed to be contiguous, then all records can be retrieved using offset computation on physical RIDs and the B-Tree is no longer necessary [21]. For variable length attributes, a standard slot array is necessary, but the attribute's position in the slot array can be used to calculate the attribute's ID. Consequently, for fixed-length attributes, the B-Tree leaf pages contain only attribute values without IDs or slot arrays, raising the effective space utilization to essentially 100%. The upper levels of the B-Tree are organized in a normal fashion with the key entry for a leaf page containing the ID value corresponding to the "smallest" tuple on the leaf page.

Processing the attribute values in a DSM partition happens in one of two ways. For a sequential scan of all values, the index is first traversed to the left-most leaf and then the leaf pages are scanned sequentially. To retrieve the attribute value for the tuple with a particular ID, the B-Tree is traversed to locate the correct leaf page by searching for the index entry that "covers" the ID. An index entry is said to cover a particular ID if the ID lies

<sup>3</sup> Handling deletes is discussed in a later section

between the index entry and its succeeding entry in the index. Once the correct leaf page is located, the page is read and the offset computation described above is used to locate the desired attribute value.

## 2.3 Tuple Reconstruction Algorithms for DSM

Scan is a fundamental database operation that scans all tuples in a table, possibly applying one or more predicates in the process. When one or more of a table's attributes are not required by subsequent operators in the query, the scan is normally combined with a project operator to eliminate unwanted attributes as output tuples are produced. In a database system that uses the NSM (or PAX) storage representation, the scan operation is trivial to implement; successive pages of a relation are read until the end of file is reached. In the case of the B-Tree DSM representation described in Section 2.2, several different scan algorithms are possible. In this section we describe and compare these algorithms for reconstructing a tuple (or portions of a tuple) from the B-Trees used to hold the vertical partitions.

### Tuple-at-a-time Reconstruction

The simplest DSM reconstruction algorithm begins by opening a sequential scan on the B-Trees of each attribute required to produce the output relation plus those attributes on which a predicate is to be applied. The scans are processed in lock-step one tuple at a time, any applicable predicates are applied and qualifying tuples are materialized in their NSM representation on the reconstruction operator's output stream. The primary disadvantage of this approach is that it incurs a random seek each time a new B-tree leaf page is read.

### Chunk-based Reconstruction

A scan of a relation stored in the DSM representation can also be viewed as a multi-way join of each of the table's vertical partitions. Since today's database systems include very efficient join algorithms, one might be tempted to simply use the standard join code to reconstruct a table from its partitions. However, reassembling a table of 20 attributes with a 19-way join is likely to overwhelm any database system. The join of the DSM partitions is actually a very special kind of merge-join in which the input tables are already sorted on the join attribute (i.e. the "virtual" ID value) and each attribute value joins exactly one attribute from all the other partitions and, thus, is handled exactly once.

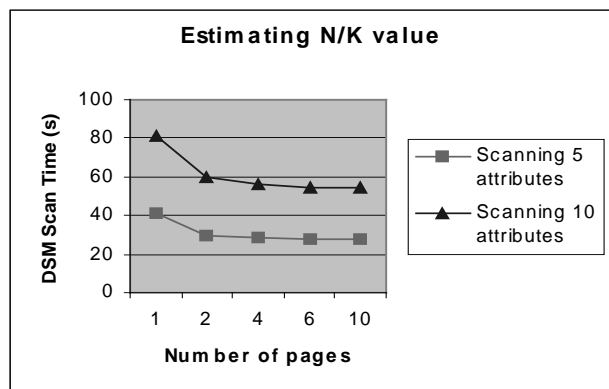
If N pages of memory are available and K attributes are accessed by the scan, the reconstruct-in-chunks algorithm begins by dividing the memory into K chunks of size N/K pages. Each chunk corresponds to one attribute. It then opens scans on the B-Trees of each of the K attributes, filling each of the K chunks with N/K leaf pages from the

corresponding B-Tree before proceeding to the next chunk. The value of this simple tactic cannot be overemphasised. While disk capacity increased a factor of 1000 in the twenty-year period from 1980 to 2000 (80 MB to 80 GB), the time for a random disk seek has decreased by only a factor of 6 (from 30 ms to 4.9 ms) over the same period. Filling each vertical partition a chunk at a time, reduces the number of random seeks performed by the “join” by a factor of N/K.

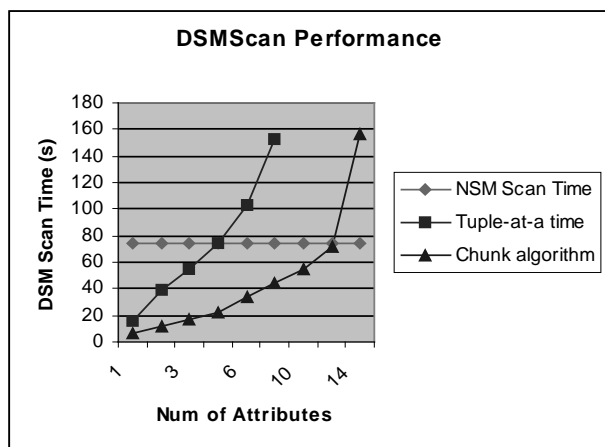
The other key technique the algorithm uses is to process attribute values in a chunk in cache-line size units to insure that the L1 and L2 data caches are used as effectively as possible. Thus, with a 64 byte cache line and 4 byte attribute values, the algorithm constructs output tuples 16 at a time.

### 2.4 Performance

This section evaluates the effectiveness of the suggested storage schemes and scan-algorithms. We first show how to select an appropriate value for N/K - the number of pages to use for each attribute with the Chunk based merge algorithm.



The graph above illustrates how the reconstruction time using the Chunk-based merge algorithm varies as a function of the chunk size for scanning 5 and 10 attributes from the 1GB version of the TPC-H Line-item table. As the graph illustrates beyond about 6 pages no further improvement occurs. For all subsequent experiments, a chunk size of 5 pages is used.



The next graph shows the DSM scan times as a function of the number of attributes being reassembled using the tuple-at-a-time and the chunk-based merge algorithms. For reference, the NSM scan time for the table is also shown.

While the tuple-at-a-time algorithm can reassemble only 4 attributes in less time than the time required to sequentially scan the entire NSM table, the chunk-based algorithm can reassemble 12 out of 16 attributes before its performance becomes worse. The results for both algorithms are much better than the results presented for the Naïve DSM implementation in Section 2.1, which required 138 seconds to reassemble just two attributes. Since the naïve representation also used a tuple-at-a-time algorithm, the primary difference is due to the improved B-tree-based storage scheme described in Section 2.2.

We think these results are very encouraging. By eliminating the redundant storage of IDs and, by using better scan algorithms, these results indicate that the DSM representation, when implemented properly, can provide better performance over a much wider range of situations than previously believed. In the following section we describe a new mirroring strategy that incorporates both NSM and DSM copies of a table.

## 3. Mirroring using DSM

### 3.1 Introduction

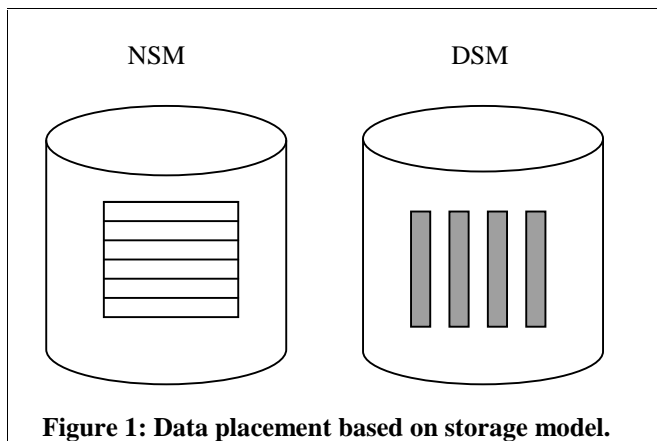
Both the NSM and DSM storage models have inherent limitations. Database systems, having to pick one, have traditionally chosen NSM, as it is more suitable for OLTP-like applications. Most database systems today employ some form of redundant storage to provide tolerance to disk failures. While RAID-5 is frequently used today, the trend is toward increased use of RAID-1 (mirroring). Even though mirroring incurs a 100% storage penalty, write operations are more efficient than RAID-5 since there is no check sum block to be updated.

In this section we describe a new form of mirroring that we term *fractured mirrors*. The basic idea is simple: rather than the two disks in a mirror being identical physically, they are instead logically identical. In particular, with fractured mirrors, one copy of each table is stored in a NSM representation and one is stored in a DSM representation. This section outlines how such a system can be constructed while retaining the advantages of both formats without losing the advantages of mirroring.

### 3.2 Data placement for Fractured Mirrors

The simplest way of implementing mirrors would be to put the NSM copy on one disk of the mirrored pair and the DSM copy on the other disk as shown in Figure 1. For

each query, the optimizer would decide which copy is best and the corresponding representation would be used to execute the query. The main disadvantage of this approach is that if the query workload is skewed towards one of the two representations, the two disks will not be utilized uniformly. Another problem is that random seeks cannot be distributed between the mirrors. This is because NSM and DSM do not have similar performance when it comes to index lookups. NSM can retrieve the entire tuple in one access, while DSM must retrieve the additional attributes by means of additional index lookups using the ID. Hence the load on the two disks will not be symmetric. It is, however, possible to place each storage model on hardware specifically tuned for the model. This is an idea to explore in the future.



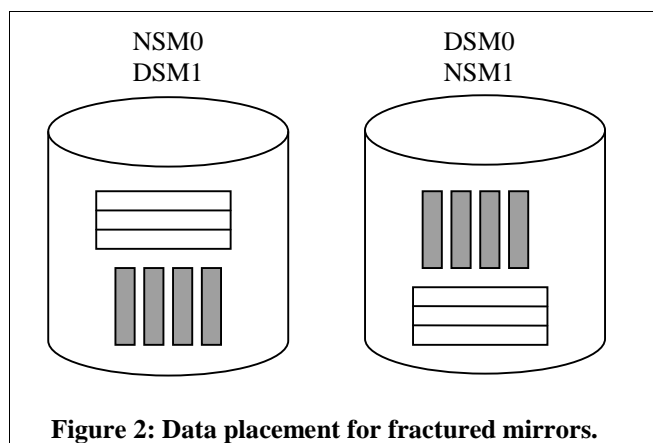
**Figure 1: Data placement based on storage model.**

A solution to this problem is the notion of *fractured mirrors*, in which data is placed on the mirrors in the following fashion. Consider a system with two disks. As shown in Figure 2, the NSM copy is declustered across the two disks using a round-robin based scheme into two equal sized fragments  $NSM_0$  and  $NSM_1$ . On disk 1, along with  $NSM_0$ , we store the tuples of  $NSM_1$  in DSM format and along with  $NSM_1$  on disk 2, we store the tuples of  $NSM_0$  in DSM format. Since both disks have  $NSM_0$  and  $NSM_1$  in some data format, they both have a complete copy of the data. Hence this constitutes a valid mirroring scheme. Even if the query workload is skewed towards one representation, since both storage formats are represented on each disk, accesses will be uniformly distributed across both disks. More importantly, we can now partition random seeks between disks in a symmetric fashion. Since the NSM copy is declustered, on average, one half of the random page accesses will be handled by each disk, a key property that the original mirroring scheme guarantees [5].

An important issue is the choice of an appropriate de-clustering algorithm [13]. It is essential that the tuples be distributed uniformly using round-robin de-clustering between the two disks, and not using a deterministic

scheme like hashing or range-partitioning. In some ways fractured mirrors are similar to RAID 10, which first mirrors an entire file and then declusters blocks between mirrors for higher bandwidth. The significant difference, of course being the presence of multiple storage representations. Another fundamental difference between the proposed system and RAID schemes is that RAID schemes usually are implemented by the disk controller in hardware. Fractured mirrors have to be implemented in software, which may lead to some inefficiency.

Given a query the database system can now select the storage format most appropriate for evaluating the query. Issues in generating query plans for the mirrors are discussed in Section 4. In the following section we present some experimental results executing queries from the TPC-H suite on this system.



**Figure 2: Data placement for fractured mirrors.**

### 3.3 Experiments on the TPC\_H Suite

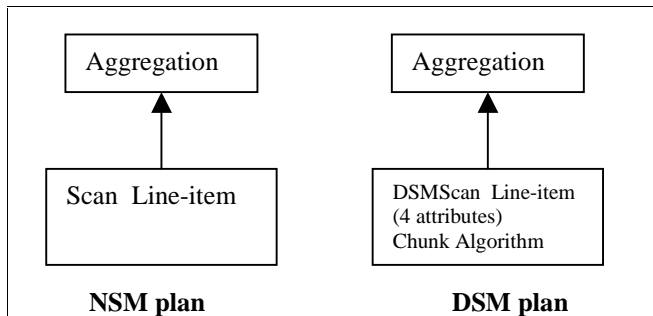
A prototype relational system was built using Shore [6] as the underlying storage manager and included the normal relational operators such as scan, join, split, merge etc along with operators to implement functionality for the chunk algorithm. The experiments were run on a Pentium III dual processor machine (550 MHz) with 1 GB of main memory running Linux 7.1. Three disks (sequential bandwidth 15-20 Mbytes/s) were used for storing data: two Shore volumes were stored on the first two for the fractured mirrors, and the third disk was used to hold the Shore log file. The Shore buffer pool size was set to 128 MB. A page size of 32 KB was used. 1 GB of TPC-H data was generated using the data generator. This data was converted into a tuple representation and stored on the two volumes as shown in Figure 2. The queries were run and their results were validated as indicated in the benchmark specification [24]. All reported times are cold times and are the average of three runs. The Shore buffer pool was flushed between queries by dismounting and remounting the disks between runs. All running times are

reported in seconds. A brief description of each query along with its execution times is given. Query plans are illustrated wherever appropriate.<sup>4</sup>

The initial queries demonstrate the advantage of maintaining a copy of the database in the DSM form. In each of these queries, the DSM plan assembles all the required partitions of a relation in a leaf node of the query plan using the Chunk Algorithm.

**Query 6:**

Query 6 computes an aggregate over selected rows of the Line-item Table. The DSM plan only scans the relevant attributes. (Only 4 attributes are used by the query)



Execution Times (s):  
 DSM: 25.555  
 NSM: 75.622

**Query 1:**

Query 1 is similar to Query 6, except it contains more complicated aggregate computations. The query touches seven of the attributes from the Line-item table and has a predicate on the l\_shipdate field that selects about 97% of the rows.

Execution Times (s):  
 DSM: 70.091  
 NSM: 143.675

**Query 12:**

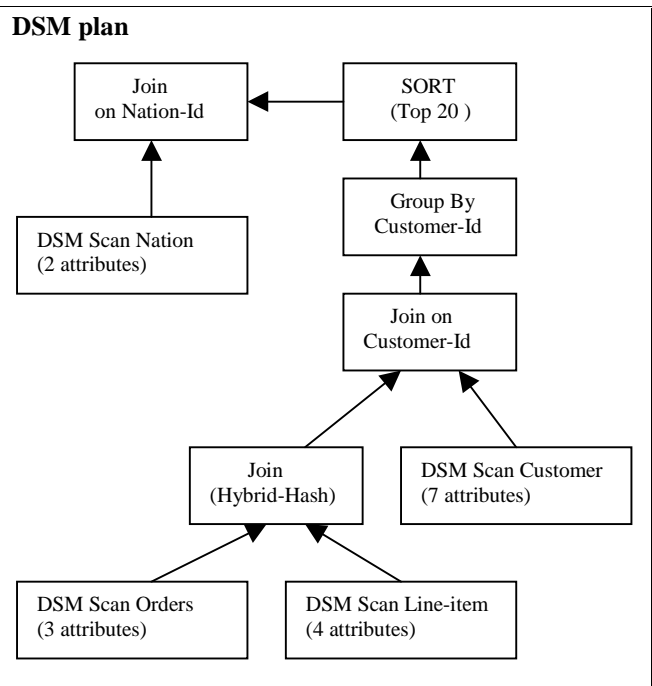
Query 12 is a join query between Line-item and Orders tables followed by an aggregation. The DSM plan consists of two DSM\_Scan nodes feeding into the join operator. Four attributes are used from Line-item and two are used from Orders.

<sup>4</sup> For simplicity, sequential plans are shown. The actual plans executed are the parallel versions taking into account the declustering in the two-disk system.

Execution Times (s):  
 DSM: 164.726  
 NSM: 232.865

**Query 10:**

Query 10 is a four-way join between the Line-item, Orders, Customer, and Nation tables. The query includes an order-by and a group-by clause and requires only the first 20 results. Order-by was implemented using the sort routine of Shore. The DSM plan is shown below with the number of attributes required from each relation. Again DSM has the best performance even though the query touches most of the attributes of the Customer table (7 out of 9).



Execution Times (s):  
 DSM: 277.416  
 NSM: 412.612

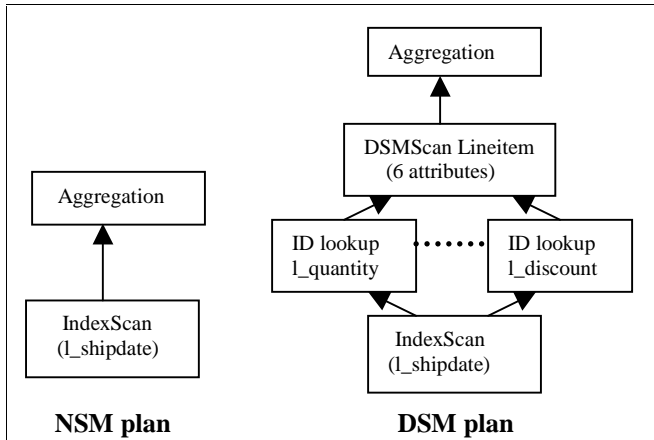
**Query 1\*:**

This query demonstrates the advantage of having both NSM and DSM representations in a mirrored system. This query is a slightly modified version of query 1 that was discussed previously. The query predicate on l\_shipdate is reversed to make it highly selective.

An index was built on the l\_shipdate column to evaluate this query efficiently. This query shows how DSM performance deteriorates with highly selective queries with even moderate degrees of projectivity since it has to probe additional indexes to fetch the required attributes.

In the DSM plan, the index on l\_shipdate produces the IDs of the qualifying tuples and then the required attributes are obtained by using the ID to probe the corresponding sparse B-Tree indexes. The 6 attributes probed are assembled using a DSMScan and the aggregate is evaluated. In this case choosing the NSM plan would be better and would be feasible in the mirrored architecture.

Execution Times (s):  
 DSM: 13.056  
 NSM: 6.455



**Query 19:**

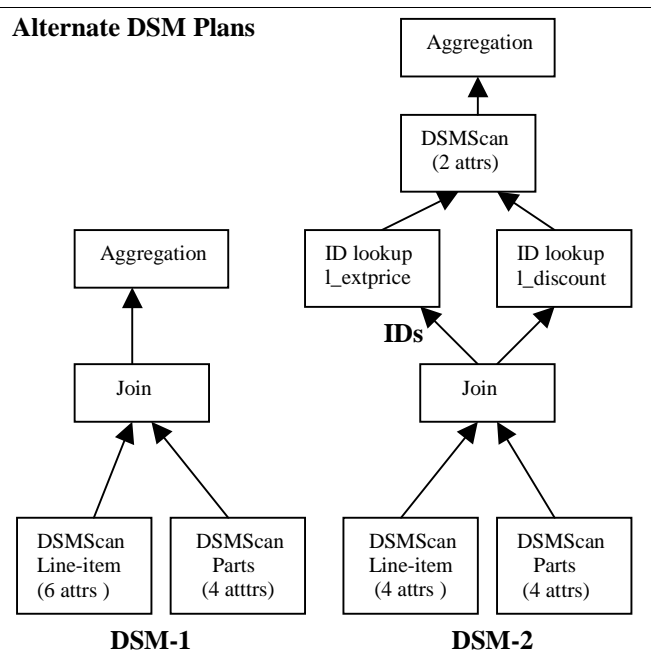
The DSM plans in the previous cases assembled all the attributes of a particular relation in a leaf node. However, in some cases, it may be more efficient to put together the attributes in multiple stages based on the selectivity of each of the attributes touched by the query.

Query 19 is a join between the Line-item table and Parts table. The query computes the revenue of parts by using the extended price and discount attributes of the Line-item table for those tuples that qualify the join. It turns out that the join is highly selective, producing only 121 tuples. Moreover, the attributes required for computing the aggregate are not required anywhere else in the plan. Hence, with the DSM plan, instead of scanning all six required attributes from Line-item in a leaf level operator (DSM-1), another plan would scan only four attributes at the leaf level (DSM-2). The IDs of the tuples produced by the join would then be used to probe the B-trees corresponding to the DSM partitions of the remaining two attributes that are needed to compute the aggregate. Since this algorithm will incur a large number of random accesses to DSM tuples, it is viable only for very highly selective predicates such as the one in this query (in which only 121 tuples out of 6 million satisfy the predicate). The DSMScan on Line-item would produce the tuple-ids along with the attributes. The join would

project the tuple-ids of the tuples that qualify the join, these ids would be used to probe the index on the partitions l\_extendedprice and l\_discount, and the values will be used to compute the aggregate. The two alternative DSM query plans are shown in the figure below.

The l\_shipmode attribute used by this query is a fixed length string type. Since the attribute contains only four distinct values, the string values are encoded as an integer field to exploit the fixed length optimizations for DSM suggested earlier.

Execution Times (s):  
 NSM: 264.469  
 DSM-1: 224.192  
 DSM-2: 208.630

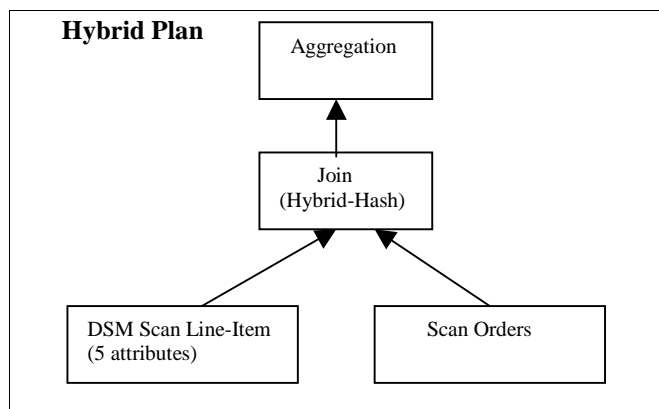


**Hybrid Plan:**

We use a simple query to demonstrate the notion of hybrid plans, plans in which both data models are used to evaluate the query plan. The query selected is a modified version of Query 12 which is a join between Line-item and Orders. An additional predicate is added to the Order table to restrict the number of order tuples and all attributes from the Orders table are projected for the tuples that qualify the join. The best plan for this query is a hybrid plan in which the NSM copy of the Orders table and the DSM copy of the Lineitem table are joined.

Execution Times (s):  
 Pure DSM: 190.804  
 Pure NSM: 193.212  
 Hybrid: 148.637

As we can see, there is not much difference between the pure NSM and pure DSM plans. Each of these plans has one leg of the join that is not optimal in terms of disk I/O. The hybrid plan uses the best means to scan each relation in the join and hence is better than the other two plans.



The speed-up obtained by using DSM for the discussed plans is summarized below. We can see that using DSM yields speed-up factors ranging from 1.3 to 3.

Query	NSM/DSM Ratio
TPC-H Query 6	2.96
TPC-H Query 1	2.04
TPC-H Query 10	1.49
TPC-H Query 1*	0.49
TPC-H Query 12	1.42
TPC-H Query 19	1.26
Hybrid Plan	1.30

We have also seen for some queries, such as Query 1\* DSM performs poorly. Having both copies as part of a mirrored system is likely to serve a wider range of query workloads. Another advantage of maintaining both representations is that the best plan for certain queries is one in which both representations are employed. It is to be noted that these numbers do not necessarily depict the best-case scenario for DSM. In an environment having relations with large number of attributes the speed-up factors could be much more substantial. For example, one of the key tables used for the Sloan Digital Sky Survey has over 400 attributes [18].

### 3.4 Synchronising the mirrors

Once the mirrors have been created, they have to be kept synchronised through the course of database operations such as inserts, updates and deletes. In traditional mirroring, all such operations are applied directly to both the copies, which is not feasible with fractured mirrors since the DSM and NSM copies do not have identical performance characteristics under these operations. For

instance, an insert operation corresponding to a tuple with  $n$  attributes would result in  $n$  insert operations on the corresponding vertical partitions of the DSM copy. Hence, given high update rates, the overhead of maintaining the mirrors up to date may result in a serious performance penalty.

The solution we are implementing uses an intermediate representation of the relation to serve as a differential file to record updates and inserts [17]. The differential file is implemented as a relation with three attributes having the schema (*Tuple-Id, Attribute-Id, Value*). A single entry represents a new attribute value of the original tuple. An insert operation would now result in the insert of  $n$  tuples to this relation, the main difference being that the inserts can be implemented as a sequence of sequential writes since the differential file is clustered on the Tuple-Id value. With main memories becoming larger and larger, the differential file can be cached in memory until the actual updates have been applied to the appropriate DSM partitions.

Once we have recorded the inserts and updates in the differential file, we have to propagate these values to the original partitions regularly to ensure that the differential file does not grow too large. Eventually we hope to piggyback these writes whenever there are reads to nearby cylinders as discussed in [16].

Deletes are handled in a slightly different fashion. We maintain a single column relation. Each page of the relation contains a bitmap. For instance a page of size 8K bytes would contain about a bitmap with 64,000 entries. To delete a particular tuple in the original relation, we need to find the page that contains the bit-entry corresponding to the given tuple-id. The index structure described in Section 2.2 can provide this access path. Once the corresponding bit has been located, it is set to 0 to indicate that the tuple has been deleted. A suitable garbage collection mechanism is used to clean up tuples that have been deleted. This is similar to the notion of an existential bitmap outlined in [21].

A side-effect of these schemes is that the differential file must be consulted during query processing. The Chunk-based reconstruction algorithm described in Section 2.3 can be extended in a simple fashion to consult the differential file and the delete bitmap while assembling tuples. The original  $k$ -way merge becomes a  $k+2$  way merge with the differential file and the delete bitmap read in tandem with the vertical partitions being assembled. Tuples from the differential file and the delete bitmap corresponding to the tuple-ids currently being reconstructed in memory are also read into main memory. We essentially ensure that the tuple being assembled has not been deleted and is also merged with the differential file updates for it before sending it to the output stream.

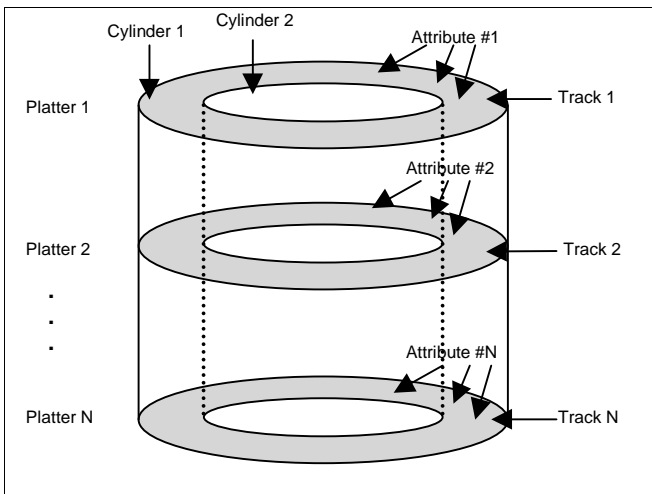


Any additionally inserted tuples in the differential file must also be processed. The differential file and the delete bitmap are clustered on tuple-id for efficient merging during the Chunk algorithm.

A disadvantage of caching the differential file in memory is that the time to reconstruct a disk after a failure may be longer than with traditional mirroring. If failure only involves a disk, the failed disk can be reconstructed using its mirror and the memory-resident differential file. If a failure involves a loss of a disk as well as the loss of memory, then it will be necessary to also use the transaction log to recovery the updates that had not yet been applied to the DSM copy on disk. The proposed scheme for handling updates should work well as long as it is possible to keep the differential file small and propagate the changes to the partitions on a regular basis.

It is possible in certain environments, such as those characterized by the TPC-C benchmark, that traditional mirroring (NSM + NSM) will likely have better performance than fractured mirrors. However, it may be the case that, if the updates are mainly updates to individual attributes, and not inserts (as is the case with TPC-C), that fractured mirrors should actually have performance similar to traditional mirrors.

Fractured mirrors are suitable for systems that have complex queries and a relatively small update to query ratio in the workload. We are currently working on an intermediate system that is likely to have update performance between the two extremes discussed. Some of the main differences of this system and the one we have discussed before are as follows. Only fixed length records are partitioned and all variable length records are clustered as a single partition. By careful data placement based on disk geometry we hope to reduce seek times.



**Figure 3: Placing partitions in adjacent cylinders**

Traditional data clustering lays out data sequentially cylinder by cylinder. We modify this slightly for the partitions. The first partition will be placed on track 1 of cylinder 1,2,3... etc. The second partition will be placed on track 2 of cylinder 1,2,3... etc. Thus, when we seek to a particular cylinder, the corresponding tracks will contain the vertical partitions of a table. With a single seek operation about 10 partitions can be reached on a modern disk drive. Given that all these are partitions contain fixed length attributes, no additional disk seeks are needed to propagate updates or inserts to these partitions. Clearly, the update performance in this approach will be intermediate to normal mirroring and the fractured mirrors approach. The data placement strategy is illustrated in Figure 3.

For this increase in performance, we need to invest more effort in data placement. If the workload characteristics are known in advance and if the update rates do not merit the increased complexity of this approach, the original scheme of fractured mirrors would be more suitable.

We intend to study compression techniques for DSM in tandem with careful data placement techniques for the partitions as future work.

#### 4. Issues in Query Processing

This section outlines how queries can be evaluated for fractured mirrors. Traditionally query processing proceeds without regard to whether the data is mirrored or not. The read requests generated during query execution are appropriately scheduled between mirrors based on expected seek times by a low-level disk scheduler. In our architecture, there is an opportunity to push this decision up to the level of the query optimizer, as it can choose a plan that better exploits the semantics of the different data formats used in the mirrors. This section explains how a traditional bottom up search based query optimizer can be extended to generate plans in this environment.

##### 4.1 Optimize-twice Approach

Since we have data stored in two different data models, a simple way to look at query optimization is to determine the best possible way to execute the plan using the NSM and DSM representations and then pick the better of the two plans. Consider two relations R (R1, R2, R3) and S (S1, S2) and a simple join query between them. (Assume R1 and S1 are the attributes on which the query is joined)

Query =  $\prod_{R2} (R \bowtie S)$  projects all R2 attributes which qualifies the join.

The corresponding DSM schema for R and S is R-1 (id, R1), R-2 (id, R2), R-3 (id, R3) and S-1 (id, S1), S-2 (id, S2)

The equivalent query for the DSM relations would be  $\prod_{R2} (R-1 \bowtie R-2 \bowtie S-1)$ .

The DSM query has the original join between the R1 partition and S1 partition and an additional join based on ID to retrieve the R2 attributes that belong to this join result. The Chunk algorithm discussed would be implemented as a specialized join algorithm that can be used for joining partitions. The query optimizer would use standard join ordering schemes and decide the best plan for the above query using a suitable cost model. Having obtained the best plans for each storage model, we can look at the optimizer cost estimates for both the plans and we would pick the plan having the better cost. Even though this approach is simple, each query is optimized twice, which would add to the overhead of query execution.

#### 4.2 Combined search of plan space

Ideally we would like to explore the search space of both storage models in a combined fashion, thereby eliminating the redundancy of the optimize-twice approach. This section outlines how a bottom-up search based optimizer can be extended to achieve this objective. A detailed overview of bottom-up search is available in the survey [10].

We need a new logical operator *Assemble* that corresponds to an operator that combines operator trees having partitions of the same relation. In the case we are assembling leaf nodes that are scans on vertical partitions, the algorithms discussed in 2.3 would be suitable for implementing this operator. In case we are assembling two trees of operators having partitions of the same original relation, the IDs of the partitions produced by the first tree would be used to probe the sparse B-Tree indexes of the partitions in the second. This operator is similar to the materialize operator proposed in [4] for evaluating pointer joins in object-oriented database systems.

Consider the simple join query considered in the previous section. First consider how the search space of plans is explored for the NSM model. For simplicity we assume that the database has no indexes to use for this query.

The given query is  $\prod_{R2} (R \bowtie S)$ . The base nodes for the search would be scan nodes on relations R and S. Let these be denoted by {R} and {S}. In the first phase of search-space exploration, all possible operators will be applied to the base set. For this example we will have a join operator that will generate the nodes {R, S} and {S, R}, corresponding to the 2 join orders possible. Since both of these nodes have the same logical properties, only the plan with the least cost will be retained. It happens that the plan chosen will be complete since it can

implement the projection on R2, which is the only operation left. Thus, it would be chosen as the optimal plan.

Let us consider combined optimization for both NSM and DSM. When starting the search, we need to add as base nodes all possible initial access paths. Hence we need to include the nodes in the previous case as well as scans on all partitions touched by the query i.e. {R}, {S}, {R1}, {R2}, {S1}. Among the set of operators that would generate new operator trees would be the join operator (as in the previous case) and the Assemble operator which would combine partitions.

The Join operator would generate joins for nodes that can satisfy the join predicate (between attributes R1 and S1). It would generate the following nodes, {R, S}, {R1, S}, {R1, S1}, {R, S1} and the corresponding nodes with the join orders reversed. The Assemble operator would generate {R1, R2}. The Assemble operator is not sensitive to the order in which the partitions are assembled. Hence {R2, R1} will not be generated. These nodes will again be grouped into equivalence classes and the minimal cost node will be retained for each class. Here are the classes and the best plans for those classes. (We do not include the nodes generated due to join commutativity for simplicity)

Class 1: {R, S} {R, S1} – best plan {R, S1}  
Class 2: {R1, S} {R1, S1} – best plan {R1, S1}  
Class 3: {R1, R2}

Among these classes, only Class 1 is complete since it can project the attribute R2 (the class includes a scan on all attributes of relation R).

In the next phase of optimization, the Assemble operator will combine {R1, S1} and {R1, S} with {R2} to generate {R1, S1, R2} and {R1, S, R2}. The join operator will combine {R1, R2} with {S} and {S1} to generate {R1, R2, S} and {R1, R2, S1}. Now, all the generated trees will belong to Class 1 and the best plan among {R, S1} (the current best plan) and these plans would be picked as the optimal plan.

The Assemble operator would maintain logical properties that combine the logical properties of the individual partitions being scanned. Among the logical properties we need to maintain is the notion of *attributes that come into scope*. For example, the join operator can evaluate the join {R1, S} because the attributes required for the join from R (just R1) have already come into scope. By maintaining this property we can ensure that vertical partitions are also considered in the joins and the Assemble operator will ensure that larger partitions are generated from smaller partitions. This ensures that the combined space of plans for NSM and DSM will be explored together.

An interesting benefit of this approach is that we can generate hybrid plans in which the final query plan has part of the query using NSM and another part using DSM. For our query example, {R, S1} is a hybrid plan. In certain cases such plans are better than fully DSM or NSM plans, which would be the only type of plans, generated in the “optimize-twice” approach.

Thus, conventional query optimization can be extended in a simple fashion to generate plans for fractured mirrors. Such an optimizer is currently being developed and we are in process of investigating further details including cost models, search space efficiency and appropriate heuristics for restricting search space.

## 5. Related Work

Disk technology trends were discussed in [9]. The authors point out that, while disk prices have dropped by a factor of 10,000, accesses per second have grown by a factor of only 100. [1] studied the importance of cache performance in query processing and PAX was proposed as a solution [2]. Given that we have an additional copy in DSM, some of the advantages of PAX can be bought by simply using the DSM copy. The effectiveness of using PAX for the NSM copy in fractured mirrors is to be studied as future work.

The notion of using DSM for good disk bandwidth and cache performance is similar to the notion of building covering indexes for the query at hand. But, for query workloads whose patterns are not known before hand, it may not be possible to build efficient covering indexes. By using the individual DSM partitions and the Chunk algorithm, we can simulate the functionality of covering indexes. Covering indexes have been studied in detail and are available in products like Microsoft SQL Server. The performance of the DSM model has been studied in [8], [12]. The conclusions were that DSM is better when the projectivity is low and the selectivity is medium to low while NSM is better when both the projectivity and selectivity are high. Performance of single attribute modification is the same for DSM and NSM, while NSM provides much better record insert/delete performance. The query processing algorithm presented in [12] needed the notion of join indexes, while we have outlined how query optimization can be extended in a general fashion to support DSM. A performance evaluation of DSM using the TPC-D benchmark has also been carried out in [22] using the Monet main memory database system.

Some of the very early prototype database systems that hinted at using decomposed storage were [19][20]. The BUBBA project was among the better-known projects that advocated the use of DSM. The BUBBA system [7] proposed a notion of using a set of inverted files and a

remainder relation as an online copy instead of mirroring. DSM has also been used as a physical storage model to implement object-oriented data models. Query rewriting schemes for translating queries on an object-based model into DSM is presented in [22] using the Monet main memory database system. The notion of projection indexes [21] is an implementation of DSM used in warehousing environments. Among today’s database products, Sybase-IQ [23] whose target market is data warehousing uses vertically partitioned attributes as its storage model. By using efficient compression techniques and advanced bitmap indexing, aggregate queries (which are typical in a warehousing environment) can be answered very efficiently. We were not able to obtain further details as to how the partitions were indexed and how queries were optimized. The fact that DSM is suitable for decision support workloads has already been discussed [21][22][23]. As far as we can tell, this paper is the first to propose the notion of mirroring using different data storage formats.

Using mirrors to optimize reads by distributing random seeks between the disks was first discussed in [5]. This technique was extended to optimize write performance. The scheme described in [14] used a notion of distorted mirrors, which worked at the granularity of disk blocks and cleverly managed the blocks in two partitions. The notion of the mirrors not being identical is similar to our general idea though their paper is not concerned with storage models. It would be interesting to see if some of their optimizations for placing disk blocks would still be valid under the current scheme of fractured mirrors. Our proposal of propagating updates to the vertical partitions by using a differential file in memory is similar to update piggybacking suggested in [16].

The 3-column relation proposed for organizing the differential file has been proposed in a different context. New e-commerce applications require data schemas that are constantly evolving and hence require table structures that are more flexible than the standard NSM representation. [3] proposed the 3-column relation as the standard storage format, whereas we use it only as a differential file to record the updates.

## 6. Conclusions and Future Work

The Decomposition Storage Model (DSM) has not found acceptance by the database vendors. Given technology trends and the need for storage architectures that are more aware of disk-arm and cache effects during query processing, DSM is likely to play an important role in the future. The paper identified some of the fundamental performance limitations of DSM. Contributions of this paper include alternate storage schemes and scan algorithms for DSM, which as demonstrated, provide a

dramatic increase in performance over the naïve implementation.

A new mirroring technique was proposed as a storage architecture that can best exploit the advantages of DSM. We would like to think of our work as extending the current spectrum of mirroring techniques. Based on the workload mix (queries and updates), the complexity of queries, and the update frequency, one can pick the mirrored architecture that is most suitable. As shown for complex queries such as those in the TPC-H suite, there are obvious benefits in maintaining a copy in DSM. For workloads that do not have high update rates, the notion of fractured mirrors is likely to suffice. For higher update rates and TPC-H like queries, the optimized version of fractured mirrors which pays more attention to data placement is likely to be a better choice. For simple queries with high update rates the original mirroring scheme is the best.

As part of our future work, we intend to examine a number of issues. Given a query workload, we need to decide good data placement schemes for the partitions. The current evaluation of the system has primarily focussed on the TPC-H query workload. Future work would include experimenting with different transaction protocols to update the DSM copy efficiently and an evaluation of the system using an OLTP benchmark like TPC-C. Query optimization for the mirrors offers many problems to be studied. We also need efficient schemes to handle variable length records and NULL values. Currently the mirroring scheme is implemented in software; it would be interesting to see if RAID hardware can be leveraged to any extent.

### Acknowledgements

We would like to thank the referees for their detailed comments, the Borg team and Anastassia Ailamaki for sharing their code and Joseph Burger for many clarifications on the efficient use of Shore.

### References

- [1] A.Ailamaki, D.DeWitt, M.Hill, D.Wood. DBMS on a modern processor: Where does Time Go? Proceedings of VLDB 1999.
- [2] A.Ailamaki, D.DeWitt, M.Hill, M.Skounakis. Weaving Relations for Cache Performance. Proceedings of VLDB 2001.
- [3] R.Agrawal, A.Somani, Y.Xu. Storing and Querying of E-Commerce Data. Proceedings of VLDB 2001.
- [4] J.Blakeley, W.J.McKenna, G.Graefe. Experiences building the open OODB Query Optimizer. Proceedings of ACM SIGMOD 1993.
- [5] D.Bitton, J.Gray. Disk Shadowing. Proceedings of VLDB 1988.
- [6] Carey et al. Shoring up Persistent Applications. Proceedings of ACM SIGMOD 1994.
- [7] G.P.Copeland, W.Alexander, E.E.Boughter, T.W.Keller. Data Placement in BUBBA. Proceedings of ACM SIGMOD 1988.
- [8] G.P.Copeland, S.Khosafian. A Decomposition Storage Model. Proceedings of ACM SIGMOD 1985.
- [9] J.Gray, G.Graefe. The 5-minute Rule Revisited and Other Storage Rules of Thumb. ACM SIGMOD Record 26(4) : 63-68 (1997)
- [10] Y.Ioannidis. Query Optimization. The Computer Science and Engineering Handbook 1997: 1038-1057
- [11] K.Keeton et al. Performance Characterization of a Quad Pentium Pro CPU using OLTP workloads: Proceedings of ISCA 1998.
- [12] S.Khosafian, G.Copeland et al. A Query Processing Strategy for the Decomposed Storage Model. Proceedings of ICDE 1987.
- [13] M.Livny, S.Khosafian, H.Boral. Multi-Disk Management Algorithms. Proceedings of SIGMETRICS 1987.
- [14] C.Orji, J.Solworth. Doubly-Distorted Mirrors. Proceedings of ACM SIGMOD 1993.
- [15] R.Ramakrishnan, J.Gerhke. Database Management Systems. WCB/McGraw-Hill 2000
- [16] J.Solworth, C.Orji. Write-only Disk Caches. Proceedings of ACM SIGMOD 1990.
- [17] D.Severance, G.Lohman. Differential Files: Their application to the maintenance of large databases. ACM TODS September 1976 vol. 1 number 3.
- [18] A.Szalay et al. Designing and mining multi-terabyte astronomy archives: The Digital Sky Survey. Proceedings of ACM SIGMOD 2000.
- [19] P.J.Tittman. An Experimental database system using binary relations. IFIP Working conference on database management. 1974
- [20] S.Todd. The PeterLee Relational Test Vehicle. IBM Systems Journal vol. 15 no 4 1976
- [21] P.O.Neil, D.Quass Improved Query Performance with Variant Indexes. Proceedings of ACM SIGMOD 1995
- [22] P.Boncz, A.N. Wilschut, M.L. Kersten Flattening an Object Algebra to Provide Performance. Proceedings of ICDE 1998
- [23] Sybase IQ White papers. [www.sybase.com](http://www.sybase.com)
- [24] TPC-H Benchmark Specification [www.tpc.org](http://www.tpc.org)