

Dynamic Maintenance of Data Distribution for Selectivity Estimation

Kyu-Young Whang, Sang-Wook Kim, and Gio Wiederhold

Received October 14, 1991; revised version received May 6, 1992; accepted May 27, 1993.

Abstract. We propose a new dynamic method for multidimensional selectivity estimation for range queries that works accurately independent of data distribution. Good estimation of selectivity is important for query optimization and physical database design. Our method employs the multilevel grid file (MLGF) for accurate estimation of multidimensional data distribution. The MLGF is a dynamic, hierarchical, balanced, multidimensional file structure that gracefully adapts to nonuniform and correlated distributions. We show that the MLGF directory naturally represents a multidimensional data distribution. We then extend it for further refinement and present the selectivity estimation method based on the MLGF. Extensive experiments have been performed to test the accuracy of selectivity estimation. The results show that estimation errors are very small independent of distributions, even with correlated and/or highly skewed ones. Finally, we analyze the cause of errors in estimation and investigate the effects of various parameters on the accuracy of estimation.

Key Words. Query optimization, physical database design, multidimensional file structure, multilevel grid files.

1. Introduction

Accurate estimation of selectivity is crucial in database query optimization and physical database design (Selinger et al., 1979; Christodoulakis, 1983; Whang and Krishnamurthy, 1990). Selectivity is defined as the ratio of the number of records that satisfy the query to the total number of records in a file. In query optimization the cost of an access plan for processing a query is estimated based on the numbers of records to be retrieved for intermediate and final results. Selectivity is used to estimate

Kyu-Young Whang, Ph.D., is Associate Professor, Computer Science Department, and Sang-Wook Kim is Director, Center for Artificial Intelligence Research, Korea Advanced Institute of Science and Technology, 373-1 Koo-Sung Dong, Yoo-Sung Ku, Daejeon, Korea; Gio Wiederhold, Ph.D., is Professor, Computer Science Department, Stanford University, Stanford, California 94305.

the number of records to be retrieved. Similarly, selectivity is used for physical database design. For example, to optimize the following query, we need accurate estimation of selectivity of the predicate $30 < \text{Age} < 35$ AND $50000 < \text{Salary} < 100000$.

```
SELECT  Name
FROM    Employee
WHERE   30 < Age < 35 AND
        50000 < Salary < 100000
```

Several selectivity estimation methods have been reported in the literature. The earliest and simplest one is based on the *uniform distribution assumption* and the *independence assumption* (Selinger et al., 1979). The former assumes that records are distributed uniformly over the domain of an attribute. The latter assumes that the distributions of different attributes are not correlated. However, these assumptions rarely hold in practical situations and thus cause significant errors in selectivity estimation. The errors are prominent, especially when the attributes are correlated (Christodoulakis, 1983; Vander Zander et al., 1986). Christodoulakis (1983) analytically showed how much the selectivity estimated under the uniform distribution and independence assumptions can deviate from the true selectivity using various parametric distributions.

To alleviate the problem caused by these simplistic assumptions, several methods based on histograms have been proposed: the equi-width method (Piatetsky and Connell, 1984), the equi-depth method (Muralikrishna and DeWitt, 1988), and the homogeneity-based method (Muthuswamy and Kershberg, 1985; Chen et al., 1990). The basic idea of the histogram approach is to capture a data distribution by dividing the domain into a set of intervals. In the equi-width method, the widths of the intervals are equal, and the number of records in each interval approximates the distribution. In the equi-depth method, the widths of the intervals are adjusted so that each interval has the same number of records. In the homogeneity-based method, the domain is partitioned into intervals in such a way that the records in each interval are uniformly distributed.

Each method has benefits and drawbacks. The equi-width method is a very simple technique corresponding to the classical histogram (Mannino et al., 1988), but has limitations. First, it is difficult to determine the boundaries of the domain and the intervals without prior knowledge of the data distribution. Second, the error in selectivity estimation can be significant, since some intervals can be heavily populated violating the assumption that the data distribution within an interval is uniform. The equi-depth method solves these problems by scanning the records of an existing file and by making each interval equal in population. However, the proposed technique cannot accommodate dynamic insertion or deletion of the records due to its static nature: the data structure has to be rebuilt periodically to accommodate changes. The homogeneity-based methods attempt to achieve high accuracy in estimation by maintaining a prespecified level of homogeneity (i.e., uniformity) within an interval. Thus, the assumption that the records are uniformly

distributed within each interval does not incur as much error in estimation as in the equi-depth method. The drawback is the static nature of the techniques requiring periodic reconstruction of the histogram. An excellent survey of various techniques appear in Mannino et al. (1988). To avoid the complexity of these approaches, some techniques attempt to use a hybrid approach (Ioannidis and Christodoulakis, 1991). For example, IBM's DB2 can record up to ten highest-frequency attribute values for special treatment. The remaining values are assumed to have a uniform distribution (Selinger, 1991).

This article presents a new selectivity estimation method for queries in the environments where records are inserted and deleted dynamically. This is an excellent method with nonuniform, highly skewed, and/or highly correlated data distributions. We also present the results of extensive experiments testing the accuracy of the proposed method, which prove to be excellent. Here we handle only range queries, leaving the cases for exact-match queries as a further study.

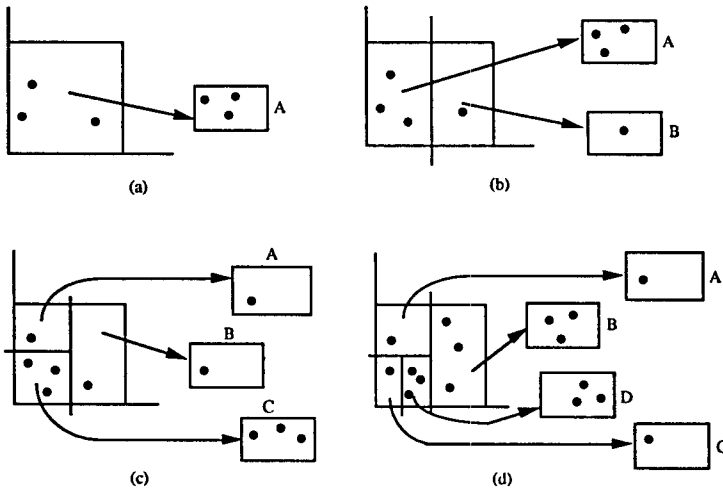
Our method uses the Multilevel Grid File (MLGF; Whang and Krishnamurthy, 1985, 1991) to maintain multidimensional data distributions. The MLGF is a dynamic, multidimensional file organization that gracefully adapts to nonuniform and correlated distributions. We discuss how each level of the MLGF directory maintains an n -dimensional data distribution, where n is the number of attributes. A lower level represents it in finer granules, and a higher level in coarser granules. Selectivity estimation is done based on the distribution information maintained in the MLGF directory.

This article is organized as follows. In Section 2 briefly reviews the structure of the MLGF. Section 3 investigates how the MLGF directory maintains a multidimensional data distribution and present the techniques for selectivity estimation. Section 4 presents experimental results obtained from testing the accuracy of estimation. Section 5 analyzes the errors in estimation and discusses how various parameters affect the errors. Section 6 briefly discusses candidate applications of the MLGF-based method. Section 7 concludes the article and proposes further study.

2. Multilevel Grid File (MLGF)

In this section we briefly review the structure of the MLGF. Section 2.1 describes the MLGF's dynamic characteristics, and Section 2.2 describes its structural characteristics. We first define some terminology: A *file* is a collection of records, where a record consists of a list of attributes. A subset of these attributes that determines the placement of the records in the file is called the *organizing attributes*. A file has a *multidimensional organization* if it contains more than one organizing attribute. A domain of an attribute is a set of values from which an attribute value can be drawn. We define the *domain space* as the Cartesian product of the domains of all the organizing attributes. We call any subset of the domain space a *region*.

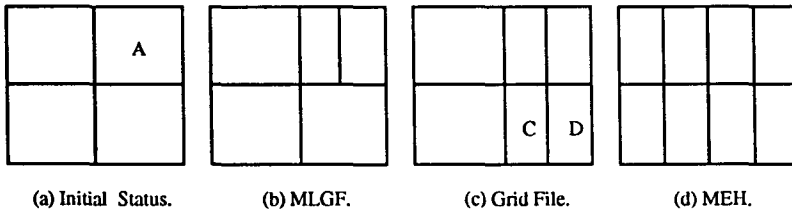
Figure 1. Dynamic growth of an MLGF



2.1 Dynamic Characteristics of MLGF

The MLGF consists of directory pages and data pages. The directory has multiple levels. An entry (*directory entry*) in the lowest level of the directory points to a data page and represents the region for which the data page is allocated. The data page contains only those records that belong to the region represented by the directory entry. The multilevel directory structure of the MLGF is built recursively; that is, a higher level of the directory is built on top of the next lower level treating it as if it were the base data.

The MLGF adapts to dynamic situations where record insertions and deletions occur by splitting and merging pages. When a new record is inserted into an n -dimensional file, the region to which the record belongs is found by searching the directory from the root to the lowest level, and the record is inserted into the data page allocated for that region. If the page overflows, the region splits into two equal-sized subregions and the records are distributed into two pages that are allocated for the new regions. Figure 1 shows how a two-dimensional MLGF grows on repeated insertion. Let us assume that a data page can contain up to three records (i.e., the data page blocking factor is three). Figure 1a is the initial state where the file contains three records in the entire domain space. If we add another record, the data page overflows causing the region to split into two, and another data page is allocated. The records are redistributed into two data pages according to the regions they belong to (Figure 1b). Figures 1c and 1d show the states of the file, when there have been subsequent splits of page A and page C.

Figure 2. Region splitting strategies

When records are deleted repeatedly, the MLGF shrinks. If the number of records in the page falls below a certain threshold (i.e., the page underflows), the region of the page is considered for merging with one of its buddies. A *buddy* of a region A is defined to be an adjacent, equal-sized region that forms a rectangular region when merged with region A. When merging actually occurs, all the records in two data pages are consolidated into one, and the other page is deallocated.

2.2 Structural Characteristics of MLGF

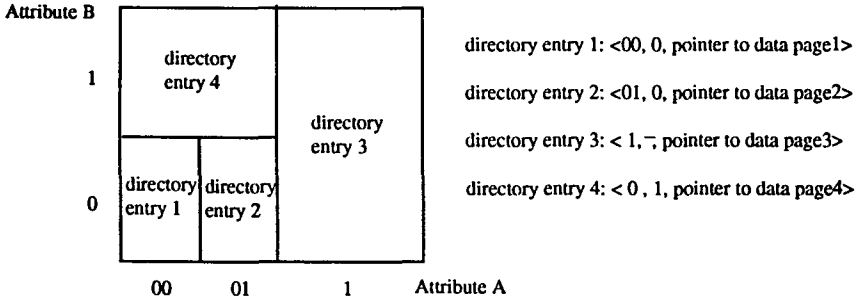
Local Splitting. The MLGF splits the region locally when the data page overflows. This is one of the advantages of using the MLGF, especially for selectivity estimation. Local splitting maintains *exactly one directory entry pointing to one data (or directory) page*. This policy makes it easier to maintain the number of records that a directory entry represents. Without it, multiple directory entries pointing to the same data or lower-level directory page have to be identified to obtain the estimation. This would be a cumbersome and costly process.

Figure 2 compares the splitting strategy of the MLGF with those of other file organizations: the grid file (Nievergelt et al., 1984) and the multidimensional extendible hashing (MEH; Otoo, 1984).

Figure 2a represents the state of domain space partition when region A is about to split. Figure 2b shows the partition that the MLGF generates after splitting region A locally. Figure 2c shows the partition that the grid file generates, in which the entire hyperplane containing region A is split. Note that, in this partition, both regions C and D point to the same data page. Figure 2d represents the partition that the MEH generates. It causes even more region splits to satisfy the equi-depth requirement of the directory for array-index computation of extendible hashing (Fagin et al., 1979). Local splitting is also employed in other data structures such as the BANG File (Freeston, 1987), the balanced multidimensional extendible hash tree (Otoo, 1986) and the K-D-B tree (Robinson, 1981).

Directory Entries. We now look into the directory structure of the MLGF in more detail. A directory entry consists of a region vector and a pointer to a data or lower-level directory page. A region vector in an n -dimensional file consists of n

Figure 3. Directory entries and the corresponding regions



hash values that uniquely identify the region. The i -th hash value of the region vector is the common prefix of the hash values for the i -th attribute of all the records that belong to the region.

For example, Figure 3 shows a partition of a two-dimensional domain space. Directory entry 1 contains the region vector $\langle 00, 0 \rangle$ that represents the lower left region and the pointer to data page 1. The hash value '00' will be the common prefix of the hash values for attribute A of the records in this region. Likewise, the hash value '0' will be the common prefix of the hash values for attribute B. Directory entry 2 contains the region vector $\langle 01, 0 \rangle$ representing a buddy of the region of directory entry 1. The symbol '-' in directory entry 3 represents the entire domain of the corresponding attribute.

A region vector also indicates the size of the region. The size of the interval in one attribute is inversely proportional to 2^v , where v is the length of the hash value of the region vector for that attribute. Therefore, the size of a region is calculated as

$$\text{region size} = \prod_{i=1}^N \frac{K(i)}{2^{v(i)}}$$

where K is the size of the entire domain of the i -th attribute.

The MLGF uses an order-preserving hashing function to map attribute values to the range of $[-2^{31}, 2^{31} - 1]$, represented by four-byte signed integers. Order-preserving hashing functions (Robinson, 1986) are generally known to be difficult to use in practice because they do not distribute values evenly over their ranges (i.e., they introduce skews in data distribution). However, data skew does not pose problems in the MLGF since its directory gracefully adapts to highly skewed or even correlated distributions. The main reason for this characteristic is the local splitting strategy of the MLGF that has exactly one directory entry to one data (or directory) page and does not represent empty regions (Whang and Krishnamurthy, 1991). The experiments in Section 4 show that selectivity estimation is not affected

significantly by data skew.

Duplicate records can be introduced due to collisions. If a data page is filled up by duplicate records, we simply chain multiple data pages, which we regard as one virtual data page. We expect, however, that hash collisions for all the organizing attributes will be infrequent.

MLGF Directory and Space Partitioning. To illustrate the hierarchical nature of the MLGF directory, consider a two-level directory consisting of D_1 and D_2 in Figure 4a. Figure 4b shows the partition of the domain space induced by each level of the directory. The round enclosures in Figure 4b are the regions represented by the directory entries in D_1 . Thus, D_1 has twelve entries. (Note that the empty region $\langle 00,1 \rangle$ was not represented in D_1 .) Level D_2 serves as the directory for level D_1 . The second directory entry in D_2 with the region vector $\langle 10,0 \rangle$ represents region b in Figure 4c and points to a page in D_1 containing three entries, which form a finer partition of region b into regions E, F, and G in Figure 4b. Note that the first (second) element of the region vector $\langle 10,0 \rangle$ is the common prefix of the first (second) elements of the region vectors of the three directory entries.

Summary. We now summarize the features of the MLGF that are relevant for selectivity estimation. Details of the MLGF operations can be found in Whang and Krishnamurthy (1985, 1991):

1. Dynamic, multidimensional file structure that adapts to nonuniform and correlated distributions.
2. Hierarchical directory structure, each level of which represents a partition of the entire domain space. The partition allows regions of different sizes.
3. Exactly one directory entry points to one data (or directory) page.
4. Empty regions not represented. This is an important property in estimating selectivity since an empty region, when merged with another region, introduces skews in data distribution within the merged region. The skew within a region is a major cause of error in estimation as we explain in Section 3.

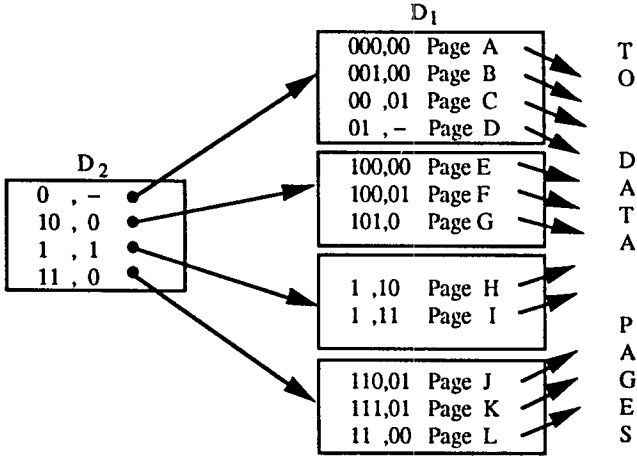
3. Selectivity Estimation Using Multilevel Grid File

In this section we first show how the MLGF directory maintains data distribution in the domain space. Next, we present the selectivity estimation method.

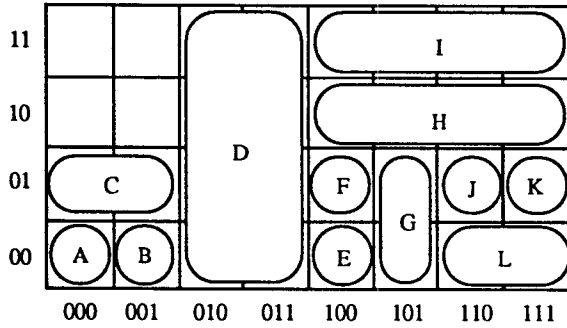
3.1 Dynamic Maintenance of Data Distribution

In Figure 4b a rounded enclosure is the region represented by a directory entry in D_1 . Because a directory entry in D_1 points to exactly one data page, the region

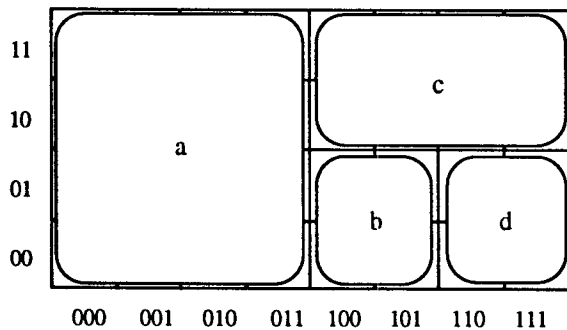
Figure 4. Two-level MLGF directory and its domain space partition



(a) The structure of a two-level MLGF directory.



(b) Regions represented by directory entries in D_1 .



(c) Regions represented by directory entries in D_2 .

corresponds to exactly one data page. We call it the *region of the data page*. Suppose that each data page contains an equal number of records. Then, each entry in the lowest level of the directory (and its region) represents an equal number of records. Thus, we derive that the density of record population is inversely proportional to the size of the region. For example, the density in region E is eight times that in region D. Therefore, the partition of the domain space derived by the lowest level of the directory (D_1) represents the data distribution. This concept is similar to that of the equi-depth histogram partitioning the domain space into buckets having the same number of records (Muralikrishna and DeWitt, 1988).

We now make a similar observation with D_2 . Assuming that each directory page contains an equal number of directory entries, we derive that each directory entry at a specific level represents an equal number of directory entries in the next lower level, and eventually represents an equal number of records. From this observation, we derive that at any specific level of the MLGF directory, the data distribution is inversely proportional to the size of the region that each directory entry represents.

In general, each level of the MLGF directory reflects the data distribution over the entire domain space. However, a lower level keeps the distribution in finer granules than a higher level does because it contains more directory entries. Thus, a lower level of the directory provides more accurate selectivity estimation. However, because each page does not necessarily contain the same number of records or directory entries, some mis-estimation is likely.

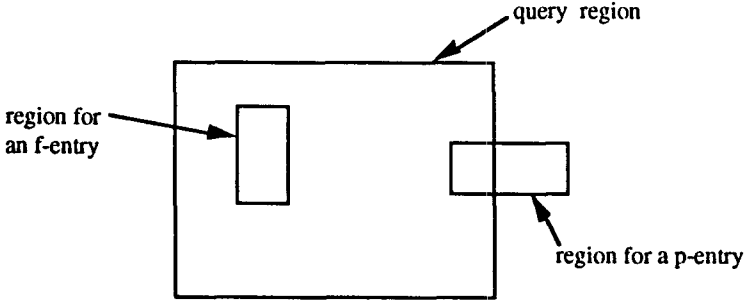
Extension of MLGF with COUNT Fields. Our model of data distribution becomes more accurate if we maintain the count of the records that each directory entry represents. One count field is kept with each directory entry. The count for a directory entry is easily maintained by updating it whenever a record or a directory entry is inserted/deleted into/from the page that the directory entry points to. It is also updated when the pages split or merge. No additional page accesses are required, although slightly more pages will be needed overall to accommodate the count fields. Maintaining these counts obviates the need for the earlier assumption that each data or directory page contains an equal number of records or directory entries. Thus, with this refined model, the data distribution is proportional to $\text{count}(i) / \text{region size}(i)$ for a directory entry i .

3.2 Selectivity Estimation

In Section 3.1 we showed how a multidimensional data distribution is derived from the MLGF directory. In this subsection we present a method for estimating selectivity from the data distribution thus obtained.

We first define some terminology. A *query region* is a subset of the domain space that satisfies the conditions of a query. We define a *full directory entry (f-entry)* as a directory entry whose region is fully enclosed by the query region, and a *partial directory entry (p-entry)* as a directory entry whose region partially overlaps the query region (Muralikrishna and DeWitt, 1988). Figure 5 shows the relationships of

Figure 5. An f-entry and p-entry versus a query region



f -entries and p -entries with the query region.

We now make the following assumption:

Assumption 1. Records are uniformly distributed within the region represented by a directory entry.

From Assumption 1 we estimate the selectivity of a query as follows:

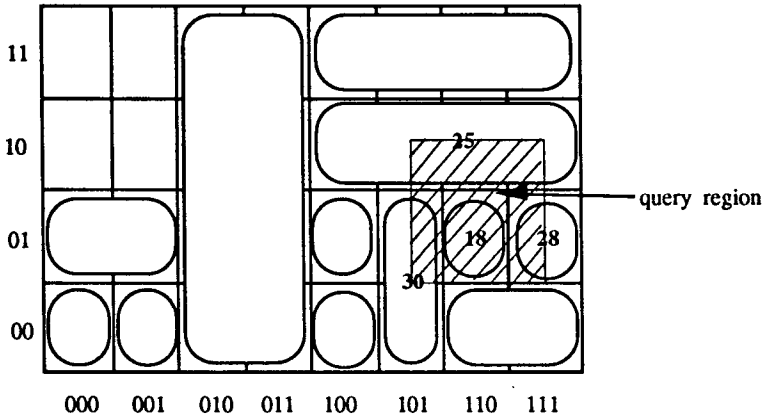
$$\text{Selectivity}(\text{query}) = \frac{\sum_{i=1}^f \text{count}(i) + \sum_{j=1}^p (\text{count}(j) \times \text{fraction}(j))}{N},$$

where f is the number of f -entries, p the number of p -entries, and N the total number of records in the file. $\text{Count}(i)$ and $\text{count}(j)$ are the record counts maintained in the directory entries i and j , respectively. The function $\text{fraction}(i)$ is defined as

$$\text{fraction}(i) = \frac{\text{size}(\text{query region} \cap \text{region of } i\text{-th } p\text{-entry})}{\text{size}(\text{region of the } i\text{-th } p\text{-entry})},$$

i.e., it is the portion of the region represented by the i -th p -entry that overlaps with the query region. The errors in estimation due to Assumption 1 are discussed in Section 5.

Example 1: Figure 6 shows the same partition of the domain space as in Figure 4b, where there are twelve directory entries. The numbers in the rounded enclosures are the counts associated with the corresponding directory entries. The file contains 300 records. Let the rectangle with slanted lines denote the query region. It contains one f -entry and three p -entries. The fractions of the p -entries overlapping with the query region are $1/4$, $1/4$ and $1/2$ (clockwise starting from the left). Thus, the estimated selectivity is $(18 + 30/4 + 25/4 + 28/2) / 300 = 0.15$.

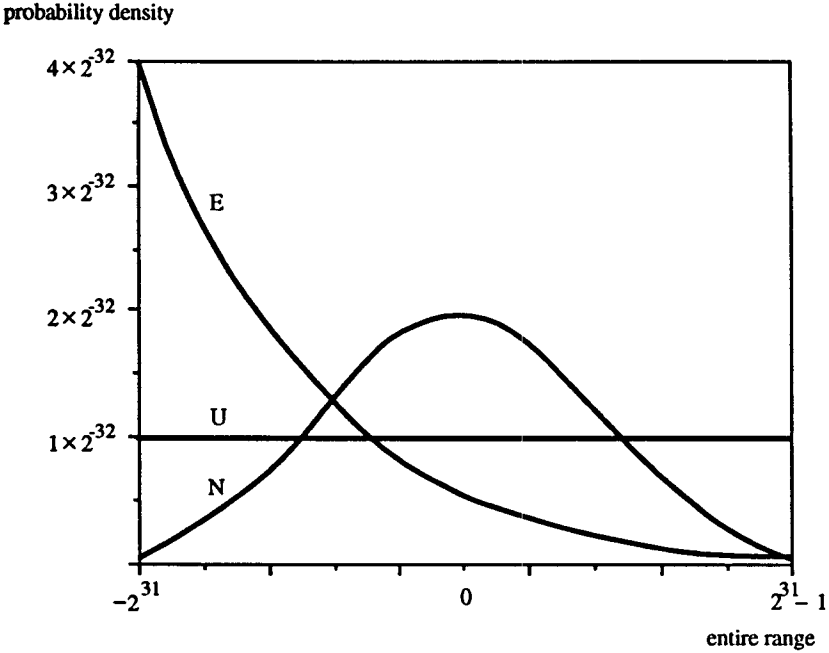
Figure 6. Query region in the domain space

4. Experimental Results

In this section we present the experimental results. The purpose of the experiments was to measure the accuracy of the selectivity estimation method under various data distributions. We also compare the results with those obtained under the uniform distribution and independence assumptions.

Two MLGF's containing 50,000 records were built: one with two organizing attributes and the other with three. The data page blocking factor (i.e., the number of records in a page; Wiederhold, 1983) was 31 maximum with an average of 21.91. The directory page blocking factor was 21 for the three-attribute file (28 for the two-attribute one) with an average of 14.81 (17.50). Each MLGF directory had three levels, with the top level 49% (33%) full.

We generated records using three categories of distributions: (1) uncorrelated distributions, (2) correlated distributions, and (3) extremely skewed distributions. Uncorrelated distributions were used to test the effect of different distributions and to help analyze the errors in estimation. Three basic one-dimensional distributions were used over the range of $[-2^{31}, 2^{31}-1]$: (1) a uniform distribution, (2) a normal distribution $N(0, \sigma^2)$, where $\sigma = 2/5 \times 2^{31}$, (3) an exponential distribution $1/\theta \times e^{-(x+2^{31})/\theta}$, where $\theta = 1/4 \times 2^{32}$. These distributions are plotted in Figure 7. The basic distributions were then composed into six combinations for the two-attribute file: UU, NN, EE, UE, UN, NE, where the first symbol indicates the distribution of the first attribute, and the second symbol of the second attribute. The symbols U, N, and E represent the uniform, normal, and exponential distribution, respectively. For the three-attribute file, there are ten such combinations. We use only three combinations, UUU, NNN, and EEE, because they produce results representative of those of the remaining combinations.

Figure 7. Basic distributions used in the experiments

A highly correlated distribution was used to test the effect of attribute correlation on the accuracy of selectivity estimation. The same normal distribution was used with the two attributes of a record set to the same value (correlation = 1). In this distribution, records reside only on the diagonal of the domain space, representing an extreme case of correlated data.

An extremely skewed distribution was used to test the effect of abnormal data skew on the accuracy of estimation. A normal distribution $N(0, \sigma^2)$ with $\sigma = 3000$ (1.4×10^{-6} of 2^{32} , the entire range) was used for both attributes. This distribution looked like a needle at the center of the two-dimensional space.

Four sets of queries were generated, each set representing a different range of selectivities: large ($1/5 \sim 2/5$), medium ($1/20 \sim 1/10$), small ($1/200 \sim 1/100$), and tiny ($1/1000 \sim 1/500$). For each set, 3,000 random queries were generated as follows: first, select the center of the initial square query region, so that the initial center points are distributed uniformly over the domain space, and then adjust the size of the query region so that it contains the number of records that satisfies the desired selectivity range. If the query region grows beyond the domain space boundary, it is clipped.

Finally, as the measure of accuracy, we used the relative error defined as follows:

$$\text{relative error} = \frac{|\text{estimated selectivity} - \text{true selectivity}|}{\text{true selectivity}}$$

Table 1. Relative errors in selectivity estimation using directory level 3 (lowest level) for 2-attribute file

		Distribution	EE	NN	NE	EU	NU	UU
Selectivity ranges	Large	avg	0.0008	0.0010	0.0009	0.0010	0.0009	0.0010
		max	0.0053	0.0042	0.0041	0.0056	0.0040	0.0054
		std	0.0006	0.0007	0.0007	0.0009	0.0007	0.0008
	Medium	avg	0.0035	0.0036	0.0032	0.0033	0.0031	0.0036
		max	0.0157	0.0180	0.0147	0.0160	0.0157	0.0186
		std	0.0027	0.0027	0.0025	0.0026	0.0025	0.0028
	Small	avg	0.0194	0.0185	0.0180	0.0178	0.0168	0.0198
		max	0.0853	0.0875	0.0729	0.0709	0.0727	0.0982
		std	0.0148	0.0142	0.0140	0.0133	0.0129	0.0159
	Tiny	avg	0.0567	0.0558	0.0573	0.0535	0.0550	0.0583
		max	0.3482	0.2908	0.2389	0.2572	0.3016	0.3209
		std	0.0453	0.0441	0.0430	0.0412	0.0440	0.0461

Table 1 shows the relative errors for the two-attribute file for six combinations of uncorrelated distributions, measured with the lowest level (Level 3) of the directory. The symbols, avg, max, std indicate the average error, maximum error, and standard deviation over 3,000 queries in each set. Because the lowest level provides the finest granularity, Table 1 shows smaller relative errors than Tables 2 and 3, where higher levels of the directory are used. We note that the average errors in Table 1 are well within an acceptable range in practice ($< 6\%$ for tiny queries) regardless of the distribution. Small standard deviations ($< 5\%$) indicate that the numbers are quite reliable.

Table 2. Relative errors in selectivity estimation using directory level 2 for 2-attribute file

		Distribution	EE	NN	NE	EU	NU	UU
Selectivity ranges	Large	avg	0.0075	0.0098	0.0065	0.0106	0.0044	0.0022
		max	0.0276	0.0334	0.0332	0.0423	0.0178	0.0154
		std	0.0050	0.0075	0.0047	0.0100	0.0036	0.0018
	Medium	avg	0.0249	0.0253	0.0195	0.0239	0.0121	0.0074
		max	0.0921	0.0744	0.0951	0.1684	0.0436	0.0374
		std	0.0195	0.0160	0.0155	0.0332	0.0088	0.0059
	Small	avg	0.1197	0.1123	0.1061	0.0851	0.0643	0.0375
		max	0.6047	0.6333	0.6771	0.3723	0.3348	0.1947
		std	0.1286	0.1237	0.1182	0.0736	0.0557	0.0293
	Tiny	avg	0.2262	0.3047	0.2618	0.1627	0.1384	0.0874
		max	1.2714	2.4556	2.6290	0.7127	0.7531	0.3848
		std	0.2415	0.4032	0.3569	0.1386	0.1158	0.0669

Table 3. Relative errors in selectivity estimation using directory level 1 (root) for 2-attribute file

Distribution		EE	NN	NE	EU	NU	UU	
Selectivity ranges	Large	avg	0.1758	0.1845	0.0770	0.2695	0.0486	0.0041
		max	0.5170	0.5524	0.2368	0.5205	0.1282	0.0180
		std	0.1061	0.1411	0.0531	0.1463	0.0292	0.0036
	Medium	avg	0.5330	0.5430	0.2271	0.2876	0.1094	0.0142
		max	1.4250	1.7532	0.7488	0.8873	0.2184	0.0549
		std	0.3471	0.3800	0.1711	0.2188	0.0614	0.0109
	Small	avg	1.4368	1.1544	0.7069	0.4794	0.2735	0.0475
		max	6.3967	5.2684	3.1316	1.7325	1.1865	0.1949
		std	1.7296	1.2431	0.7420	0.3917	0.2559	0.0333
	Tiny	avg	1.9590	1.5789	1.0735	0.5555	0.4072	0.0904
		max	13.8469	9.1387	8.2103	2.5005	2.7826	0.3956
		std	2.8761	1.9617	1.5184	0.5072	0.4802	0.0716

For comparison, we show the errors when selectivity is estimated using the uniform distribution assumption (Table 4). As expected, the errors are large (up to 852% in average errors) and vary widely depending on the distribution. From Tables 1, 2, 3 we can observe that errors get larger as a higher level of the directory is used, converging to those in Table 4. This tendency is expected because a higher level of the directory represents the distribution in coarser granules, within which the distributions are assumed uniform.

Table 4. Relative errors in selectivity estimation using uniform distribution assumption for 2-attribute file

Distribution		EE	NN	NE	EU	NU	UU	
Selectivity ranges	Large	avg	1.0532	0.2133	0.6329	0.6923	0.1816	0.0062
		max	2.6177	0.7199	1.8875	1.9485	0.4351	0.0222
		std	0.6047	0.2018	0.4181	0.4465	0.1309	0.0044
	Medium	avg	3.1623	0.6692	1.4646	1.5410	0.5343	0.0169
		max	10.5084	1.4543	3.8861	4.6248	1.4146	0.0646
		std	3.0921	0.3820	1.2312	1.3858	0.3390	0.0112
	Small	avg	6.6049	1.9261	3.1990	2.2944	1.0795	0.0479
		max	77.9682	8.3329	16.5428	9.5391	4.3592	0.2197
		std	9.1933	2.1457	3.8633	2.5578	1.1722	0.0346
	Tiny	avg	8.5256	2.8499	4.4684	2.4858	1.3374	0.0905
		max	61.6664	20.7783	38.1066	12.3821	8.2856	0.4069
		std	13.2175	4.0964	6.9647	2.9640	1.7147	0.0725

Tables 5 to 7 show the results for the three-attribute file. We see the same tendencies as in Tables 1 to 3, and errors in similar ranges ($< 10.2\%$ in average errors). The errors in the three-attribute file are slightly larger than those for the two-attribute file, because a data page represents a larger interval in the domain of each attribute, thus causing more difference in distribution within a data page (Section 5). Table 8 shows the results of the uniform distribution assumption for the three-attribute file. We observe the same tendencies as in the two-attribute file. (The average error is as high as 1,381%.)

Table 5. Relative errors in selectivity estimation using directory level 3 (lowest level) for 3-attribute file

Distribution		EEE	NNN	UUU	
Selectivity ranges	Large	avg	0.0056	0.0078	0.0017
		max	0.0260	0.0230	0.0083
		std	0.0044	0.0048	0.0013
	Medium	avg	0.0156	0.0111	0.0054
		max	0.0583	0.0373	0.0244
		std	0.0117	0.0080	0.0042
	Small	avg	0.0449	0.0417	0.0253
		max	0.2011	0.2292	0.1054
		std	0.0381	0.0369	0.0192
	Tiny	avg	0.0950	0.1020	0.0737
		max	0.4790	0.5965	0.3153
		std	0.0789	0.1015	0.0578

Table 6. Relative errors in selectivity estimation using directory level 2 for 3-attribute file

Distribution		EEE	NNN	UUU	
Selectivity ranges	Large	avg	0.0372	0.0391	0.0032
		max	0.1410	0.1122	0.0119
		std	0.0291	0.0212	0.0023
	Medium	avg	0.0789	0.0537	0.0092
		max	0.3268	0.1505	0.0426
		std	0.0720	0.0319	0.0070
	Small	avg	0.2486	0.8192	0.0380
		max	0.9054	0.7831	0.1577
		std	0.2146	0.1522	0.0290
	Tiny	avg	0.4157	0.4175	0.1005
		max	1.8329	1.9872	0.4718
		std	0.3823	0.4021	0.0719

Table 7. Relative errors in selectivity estimation using directory level 1 (root) for 3-attribute file

		Distribution	EEE	NNN	UUU
Selectivity ranges	Large	avg	0.1536	0.2735	0.0048
		max	0.4375	0.7464	0.0188
		std	0.1073	0.1407	0.0033
	Medium	avg	0.4575	0.3566	0.0125
		max	1.1848	1.0818	0.0500
		std	0.2825	0.2119	0.0091
	Small	avg	1.4833	1.2404	0.0445
		max	9.5092	5.6285	0.1976
		std	1.5126	1.1822	0.0333
	Tiny	avg	3.3038	2.1784	0.1055
		max	24.5124	13.8064	0.4713
		std	4.5348	2.6542	0.0729

Table 8. Relative errors in selectivity estimation using uniform distribution assumption for 3-attribute file

		Distribution	EEE	NNN	UUU
Selectivity ranges	Large	avg	0.9631	0.3997	0.0048
		max	2.4754	0.8071	0.0202
		std	0.5723	0.1486	0.0036
	Medium	avg	3.3354	0.3880	0.0131
		max	10.1056	0.8514	0.0575
		std	2.6335	0.2085	0.0097
	Small	avg	9.0556	1.5060	0.0453
		max	63.2002	5.0010	0.1884
		std	12.0549	1.3610	0.0335
	Tiny	avg	13.8088	2.8187	0.1057
		max	209.8033	15.0495	0.4777
		std	23.4148	3.2901	0.0728

Table 9 shows the results of the correlated distribution and Table 10 of the uniform distribution assumption (for all three levels). Table 9 indicates that our method works as well with a highly correlated distribution as with a uniform distribution if the lowest level of the directory is used (see the case of UU in Table 1). The errors from the uniform distribution assumption (Table 10) are very large (up to 6,628%) in this case since the data reside only on the diagonal of the domain space.

Table 9. Relative errors in selectivity estimation using correlated distribution (NN) where attribute 1 = attribute 2

		Directory Level	Level 1	Level 2	Level 3
Selectivity ranges	Large	avg	0.1580	0.0016	0.0002
		max	0.5699	0.0256	0.0016
		std	0.1313	0.0031	0.0002
	Medium	avg	0.5288	0.0083	0.0010
		max	2.2166	0.1205	0.0062
		std	0.3888	0.0156	0.0010
	Small	avg	5.1674	0.1136	0.0097
		max	22.7169	1.1463	0.0544
		std	5.2772	0.1485	0.0095
	Tiny	avg	27.1087	1.3569	0.0509
		max	111.4295	5.6396	0.3046
		std	26.4999	0.9417	0.0510

Table 10. Relative errors in selectivity estimation using uniform distribution for correlated distribution (NN) assumption where attribute 1 = attribute 2

Selectivity ranges	Large	avg	0.3150
		max	0.9344
		std	0.2405
	Medium	avg	1.0999
		max	4.1898
		std	0.9377
	Small	avg	12.9722
		max	48.8245
		std	12.7237
	Tiny	avg	66.2761
		max	244.1935
		std	63.6555

Table 11 shows the results of testing the extremely skewed distribution for all three levels. Table 12 presents the results of selectivity estimation using the uniform distribution assumption. The errors are generally very small at level 3 except for a slightly higher error (an average error of 40.6%) for queries with tiny selectivity ranges. The error is an aberration due to the simple design of the query distribution (uniform distribution). Since the generated queries are uniformly distributed over the domain space, most queries fall in the periphery of the populated area resulting

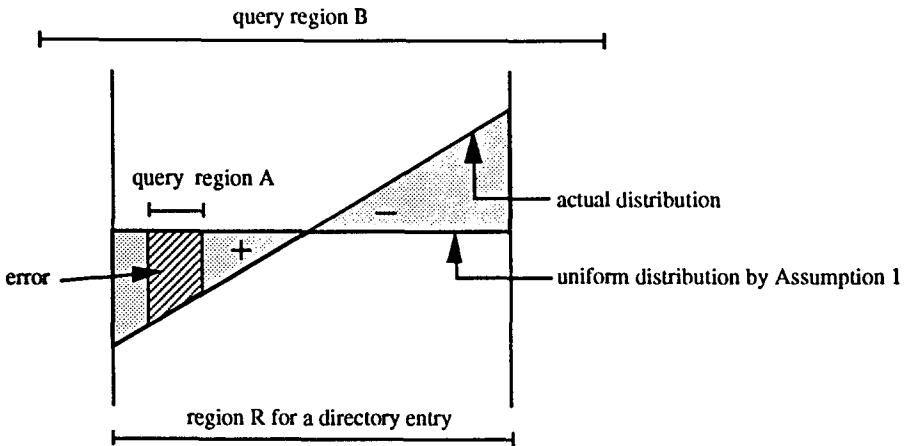
in over-emphasis on the periphery. A small number of data pages in the periphery (usually under-populated) tend to have large distribution changes, which cause errors as explained in Section 5. Overall, however, our technique provides an excellent estimation even with highly skewed distributions. In comparison, the result of the uniform distribution assumption contains an average error as high as 21,847%.

Table 11. Relative errors in selectivity estimation for extremely skewed distribution (NN) where Standard Deviation $\sigma = 3,000$

Directory Level		Level 1	Level 2	Level 3	
Selectivity ranges	Large	avg	0.4094	0.0111	0.0009
		max	0.5288	0.0153	0.0013
		std	0.1074	0.0033	0.0002
	Medium	avg	1.0934	0.0798	0.0109
		max	1.9030	0.0956	0.0153
		std	0.7599	0.0125	0.0033
	Small	avg	6.7251	0.7791	0.0640
		max	11.0899	1.2487	0.0922
		std	4.1915	0.4611	0.0273
	Tiny	avg	13.8618	2.4938	0.4062
		max	25.4795	4.2482	0.8087
		std	8.4702	1.1423	0.2248

Table 12. Relative errors in selectivity estimation using uniform distribution assumption for extremely skewed distribution (NN) where Standard Deviation $\sigma = 3,000$

Selectivity ranges	Large	avg	0.4629
		max	1.1681
		std	0.2713
	Medium	avg	3.2485
		max	6.5203
		std	1.5146
	Small	avg	37.5739
		max	70.8783
		std	15.3781
	Tiny	avg	218.4697
		max	468.4087
		std	93.6126

Figure 8. Estimation error caused by nonuniform distribution

5. Error Analysis

In this section we analyze the cause of estimation errors. Assumption 1 states that records are uniformly distributed within the region represented by a directory entry. In practice, however, this assumption is often violated, leading to errors.

Consider Figure 8, where one-dimensional query regions A and B overlap with region R for a directory entry. The actual distribution within R is nonuniform, but is assumed uniform by Assumption 1. Thus, for query region A, the area with slanted lines represents the error in estimation. For query region B, region R does not incur any error because it is completely contained in B. In this case the errors cancel each other out, because both the actual and uniform distributions represent the same number of records over the entire region R. We have the following Lemma:

Lemma 1. An f -entry does not contribute errors in selectivity estimation.

Thus, only p -entries contribute errors.

There are several parameters that affect estimation errors:

1. *Size of the query region.* Lemma 1 explains why the errors are smaller for larger queries, which have proportionally more f -entries and fewer p -entries than smaller queries. In all the tables we observe that queries with tiny selectivity ranges have large errors.
2. *Level of the directory.* Errors are smaller at a lower level of the directory, since the domain space is partitioned into regions in finer granules and thus the distribution is more uniform within a region. In the experiments, we observe that the lowest level (Level 3) produces the best results. This is hardly surprising since more information is available at low levels of the directory.

3. *Number of records in a file.* The errors are smaller when the file contains more records (i.e., when the database is large). Since a larger database contains more data pages, the regions become smaller. Thus, for the same reason as in item 1, the error is reduced.
4. *Data page blocking factor.* A smaller blocking factor reduces the error since it induces smaller regions at the lowest level of the directory.

6. Applications

In this section we discuss alternative file configurations that can be applied to practical situations. The first alternative is to use the MLGF to organize the file, as well as to estimate the selectivity. In this case, the MLGF serves as the multidimensional index structure (we need only one) that replaces multiple single- (or multiple-) column indexes. The clustering property of such an index for a single column is somewhat worse than the dedicated clustering index (Astrahan, et al., 1976), but is much better than a nonclustering index (see discussions on partial-match queries in Whang and Krishnamurthy, 1991). The reason for this is that, in the MLGF, the clustering property is shared evenly by all the organizing attributes.

The second alternative is to use a conventional index as the primary (clustering) index and to use the MLGF as a secondary index for faster access through non-indexed attributes as well as for selectivity estimation. In this case, the MLGF must be a full index (Wiederhold, 1987) (a maximum-resolution directory in Whang and Krishnamurthy, 1991) providing a directory entry pointing to each record at the lowest level. A full index is necessary because the MLGF as a secondary index cannot dictate the placement of the records. Instead, the directory entries at the lowest level serve as the surrogates of the records, whose placement is determined by the MLGF structure. In this case, the selectivity estimation is done using any level of the directory except for the lowest level. This alternative has an advantage of incorporating an MLGF on top of existing files without disrupting their structures. The overhead of maintaining an MLGF index is approximately the same as that of maintaining a B-tree index. If the MLGF contains multiple attributes, the cost is somewhat higher since the entire region vector is compared at each step of searching.

7. Conclusion

We proposed a novel approach to multidimensional selectivity estimation. The key element in this approach is dynamic maintenance of multidimensional data distribution. We showed that each level of the MLGF directory naturally reflects a multidimensional data distribution, where a lower level of the directory provides finer granularity and a higher one coarser granularity. The estimated distribution has been further refined by employing a COUNT field for each directory entry. To

the extent of the authors' knowledge, the idea of estimating selectivity based on dynamic maintenance of multidimensional data distribution is new.

We showed, through extensive experiments, that the proposed selectivity estimation method works excellently independent of distributions, even with correlated and/or highly skewed ones. Results show that estimation errors are practically small (average errors $< 10.2\%$ for most distributions and $< 40.6\%$ for the extremely skewed one) when the lowest level of the directory is used.

We analyzed the cause of errors and investigated the effects of various parameters on the accuracy of estimation. We concluded that errors decrease for a large database, a low level of the directory, a small data page blocking factor, and a large query region.

As a further study, we are considering extending this approach to exact-match queries involving equality predicates. Here, we need a technique for handling duplicates of distinct values in a discrete distribution (Whang et al., 1990). We also plan to perform a more detailed error analysis and to provide recommendations for practitioners on such parameters as page size, blocking factor, and directory level to be selected to satisfy specific requirements.

Acknowledgment

This research was partially supported by NSF Grant IRI907733 while the first two authors were visiting the Computer Science Department, Stanford University, during the summer of 1991. The work was supported by KOSEF (Korea Science and Engineering Foundation) Grant 911-1102-001-2 through the Korea Advanced Institute of Science and Technology. The authors wish to thank Yitzhak Birk for reminding them of count fields, which helped make the estimation method more accurate. Sang K. Cha carefully read the preliminary version of this article and contributed many useful comments.

References

- Astrahan, M.M., Blasgren, M.W., Chamberlin, D.D., Eswaran, K.P., Gray, J.N., Griffiths, P.P., King, W.R., Lorie, R.A., McJones, P.R., Mehl, J.W., Putzolu, G.R., Traiger, I.L., Wode, B.W., and Watson, V. System R: Relational approach to database management. *ACM Transactions on Database Systems*, 1(2):97-137, 1976.
- Chen, M.C., McNamee, L., and Matloff, N. Selectivity estimation using homogeneity measurement. *Proceedings of the International Conference on Data Engineering*, Los Angeles, 1990.
- Christodoulakis, S. Estimating record selectivities. *Information Systems*, 8(2):105-115, 1983.

- Fagin, R., Nievergelt, J., Pippenger, N., and Strong, H.R. Extendible hashing: A fast access method for dynamic files. *ACM Transactions on Database Systems*, 4(3):315-344, 1979.
- Freeston, M. The BANG file: A new kind of grid file. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, San Francisco, 1987.
- Ioannidis, Y.E. and Christodoulakis, S. On the propagation of errors in the size of join results. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, Denver, Colorado, 1991.
- Mannino, M.V., Chu, P., and Sagar, T. Statistical profile estimation in database systems. *ACM Computing Surveys*, 20(3):191-221, 1988.
- Muralikrishna, M. and DeWitt, D. Equi-depth histograms for estimating selectivity factors for multi-dimensional queries. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, Chicago, 1988.
- Muthuswamy, B. and Kershberg, L. A DDSM for relational query optimization. *Proceedings of the ACM Annual Conference*, Denver, Colorado, 1985.
- Nievergelt, J., Hinterberger, H., and Sevcik, K.C. The grid file: An adaptable, symmetric multikey file structure. *ACM Transactions on Database Systems*, 9(1):38-71, 1984.
- Otoo, E.J. A mapping function for the directory of a multidimensional extendible hashing. *Proceedings of the Tenth International Conference on Very Large Data Bases*, Singapore, 1984.
- Otoo, E.J. Balanced multidimensional extendible hash tree. *Proceedings of the ACM Symposium on Principles of Database Systems*, Cambridge, MA, 1986.
- Piatetsky, S.G. and Connell, G. Accurate estimation of the number of tuples satisfying a condition. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, New York, 1984.
- Robinson, J.T. The K-D-B tree: A search structure for large multidimensional dynamic indexes. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, New York, 1981.
- Robinson, J.T. Order preserving linear hashing using key statistics. *Proceedings of the ACM Symposium on Principles of Database Systems*, Cambridge, MA, 1986.
- Selinger, P.G., Astrahan, M.M., Chamberlin, D.D., Lorie, R.A., and Price, T.G. Access path selection in a relational database management system. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, San Jose, California, 1979.
- Selinger, P.G. private communication, 1991.
- Vander Zander, B., Taylor, H., and Bitton, D. Estimating block accesses when attributes are correlated. *Proceedings of the Twelfth International Conference on Very Large Data Bases*, Kyoto, Japan, 1986.
- Whang, K.Y. and Krishnamurthy, R. Multilevel grid files. *IBM Research Report RC 11516*, 1985.

- Whang, K.Y. and Krishnamurthy, R. Query optimization in a memory-resident domain relational calculus database systems. *ACM Transactions on Database Systems*, 15(1):67-95, 1990.
- Whang, K.Y., Vander-Zander, B., and Taylor, H. A linear-time probabilistic counting algorithm for database applications. *ACM Transactions on Database Systems*, 15(2):208-229, 1990.
- Whang, K.Y. and Krishnamurthy, R. The multilevel grid file—A dynamic hierarchical multidimensional file structure. *Proceedings of the International Conference on Database Systems for Advanced Applications*, Tokyo, Japan, 1991.
- Wiederhold, G. *Database Design*. 2nd ed. New York: McGraw-Hill Book Company, 1983.
- Wiederhold, G. *File Organization for Database Design*, New York: McGraw-Hill Book Company, 1987.