# Modeling Multidimensional Databases

*Rakesh Agrawal    Ashish Gupta*    *Sunita Sarawagi*

IBM Almaden Research Center
650 Harry Road, San Jose, CA 95120

## Abstract

*We propose a data model and a few algebraic operations that provide semantic foundation to multidimensional databases. The distinguishing feature of the proposed model is the symmetric treatment not only of all dimensions but also measures. The model provides support for multiple hierarchies along each dimension and support for adhoc aggregates. The proposed operators are composable, reorderable, and closed in application. These operators are also minimal in the sense that none can be expressed in terms of others nor can any one be dropped without sacrificing functionality. They make possible the declarative specification and optimization of multidimensional database queries that are currently specified operationally. The operators have been designed to be translated to SQL and can be implemented either on top of a relational database system or within a special purpose multidimensional database engine. In effect, they provide an algebraic application programming interface (API) that allows the separation of the frontend from the backend. Finally, the proposed model provides a framework in which to study multidimensional databases and opens several new research problems.*

## 1   Introduction

Codd [CCS93] coined the phrase On-Line Analytical Processing (OLAP) to characterize the requirements for summarizing, consolidating, viewing, applying formulae to, and synthesizing data according to multiple dimensions. OLAP software enables analysts, managers, and executives to gain insight into the performance of an enterprise through fast access to a wide variety of views of data organized to reflect the multidimensional nature of the enterprise data [Col95]. It has been said that current relational database systems have been designed and tuned for On-Line Transaction Processing (OLTP) and are in-

adequate for OLAP applications [Cod93] [Fin95] [KS94]. In response, several multidimensional database products have appeared on the market (see [Rad95] for a survey).

The database research community, however, has so far not played a major role in this market phenomenon. Gray et al. [GBLP96] recently proposed an extension to SQL with a *Data Cube* operator that generalizes the groupby construct to support some of the multidimensional analysis. Since then, techniques have been developed for computing the data cube [AAD+96], for deciding what subset of a data cube to pre-compute [HRU96] [GHRU97], for estimating the size of multidimensional aggregates [SDNR96], and for indexing pre-computed summaries [SR96] [JS96]. The research in multidimensional indexing structures (see, for example, [Gut94] for an overview) is relevant as well. Lastly, research in statistical databases (see, for example, [Sho82] for an overview) also addressed some of the same concerns.

This paper presents a framework for research in multidimensional databases. We first review concepts and terminologies in vogue in multidimensional database products in Section 2. We also point out some of the deficiencies in the current products. We then propose in Section 3 a data model to provide semantic backing to the techniques used by current multidimensional database products. The salient features of our model are:

- Our data model is a multidimensional cube with a set of basic operations designed to unify the divergent styles in use today and to extend the current functionality.

- The proposed model provides symmetric treatment to not only all dimensions but also to measures. The model also is very flexible in providing support for multiple hierarchies along each dimension and support for adhoc aggregates.

- Each of our operators are defined on the cube and produce as output a new cube. Thus the operators are closed and can be freely reordered. This free composition allows a user to form larger queries, thereby replacing the relatively inefficient one-operation-at-

a-time approach of many existing products. The algebraic nature of the cube also provides an opportunity for optimizing multidimensional queries.

- The proposed operators are minimal. None can be expressed in terms of others nor can any one be dropped without sacrificing functionality.

- Our modeling framework provides the logical separation of the frontend graphical user interface (GUI) used by a business analyst from the backend storage system used by the corporation. The operators thus provide an algebraic application programming interface (API) that allows the interchange of frontends and backends.

We discuss in Section 4 some of our design choices and show how currently popular multidimensional database operations can be expressed in terms of the proposed operators. These operators have been designed to be translated into SQL, albeit with some minor extensions. We refer the reader to [AGS96] for these translations. Thus, our data model can be implemented on either a general-purpose relational system or a specialized engine. We conclude with a summary in Section 5.

## 2 Current State of the Art

We begin with a brief overview of the current state of art in multidimensional databases.

**Example 2.1** Consider a database that contains point of sale data about the sales price of products, the date of sale, and the supplier who made the sale. The $sales$ value is functionally determined by the other three attributes. Intuitively, each of the other three attributes can "vary" and accordingly determine the sales value. Figure 1 illustrates this "multidimensional" view.
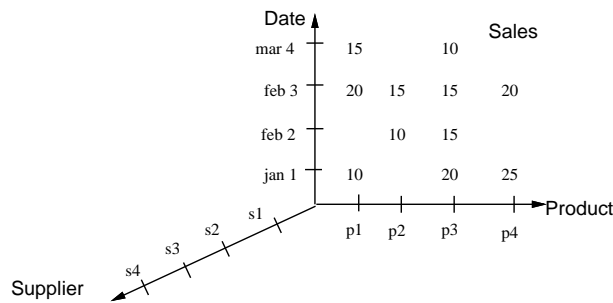


**Figure 1. Example data cube**

$\square$

### 2.1 Terminology

Determining attributes like $product$, $date$, $supplier$ are referred to as *dimensions* while the determined at-

tributes like $sales$ are referred to as *measures*. (Dimensions are called categorical attributes and measures are called numerical or summary attributes in the statistical database literature [Sho82]). There is no formal way of deciding which attributes should be made dimensions and which attributes should be made measures. It is left as a database design decision.

Dimensions usually have associated with them *hierarchies* that specify aggregation levels and hence granularity of viewing data. Thus, $day \rightarrow month \rightarrow quarter \rightarrow year$ is a hierarchy on $date$ that specifies various aggregation levels. Similarly, $product\ name \rightarrow type \rightarrow category$ is a hierarchy on the $product$ dimension.

An analyst might want to see only a subset of the data and thus might view some attributes and within each selected attribute might restrict the values of interest. In multidimensional database parlance, the operations are called *pivoting* (rotate the cube to show a particular face) and *slicing-dicing* (select some subset of the cube). The multidimensional view allows hierarchies associated with each dimension also to be viewed in a logical manner. Aggregating the product dimension from product name to product type is expressed as a *roll-up* operation. The converse of roll-up is *drill-down* that displays detail information for each aggregated point. Thus, drilling-down the product dimension from product category to product type gets the sales for each product type corresponding to each product category. Further drill down will get sales for individual products. Drill-down is essential because often users want to see only aggregated data first and selectively see more detailed data.

**Example 2.2** We give below some queries to provide a flavor of multidimensional queries. These queries use the database from Example 2.1 and other necessary hierarchies on product and time dimensions.

- Give the total sales for each product in each quarter of 1995. (Note that quarter is a function of date).

- For supplier "Ace" and for each product, give the fractional increase in the sales in January 1995 relative to the sales in January 1994.

- For each product give its market share in its category today minus its market share in its category in October 1994.

- Select top 5 suppliers for each product category for last year, based on total sales.

- For each product category, select total sales this month of the product that had highest sales in that category last month.

- Select suppliers that currently sell the highest selling product of last month.

- Select suppliers for which the total sale of every product increased in each of last 5 years.

- Select suppliers for which the total sale of every product category increased in each of last 5 years.

□

## 2.2 Implementation Architectures

There are two main approaches currently used to build multidimensional databases. One approach maintains the data as a $k$-dimensional matrix based on a non-relational specialized storage structure. The database designer specifies all the aggregations they consider useful. While building the storage structure, these aggregations associated with all possible roll-ups are precomputed and stored. Thus, roll-ups and drill-downs are answered in interactive time. Many products have adopted this approach – for instance, Arbor Essbase [Arb] and IRI Express [IRI].

Another approach uses a relational backend wherein operations on the data cube are translated to relational queries (posed in a possibly enhanced dialect of SQL). Indexes built on materialized views are heavily used in such systems. This approach also manifests itself in many products – for instance, Redbrick [Eri95] and Microstrategy [Mic].

## 2.3 Additional Desired Functionality

We believe that multidimensional database systems must provide the following additional functionality, which is either missing or poorly supported in current products:

- **Symmetric treatment not only of all dimensions but also of measures**. That is, selections and aggregations should be allowed on all dimensions and measures. For example, consider a query that finds the total sales for each product for ranges of sales price like 0-999, 1000-9999 and so on. Here the sales price of a product, besides being treated as a measure, is also the grouping attribute. Such queries that require categorizing on a "measure" are quite frequent. Non-uniform treatment of dimensions and measures makes such queries hard in current products. The proposed OLAP council API [OLA96] provides for all the measures to be put on one dimension of the cube. However, this proposal maintains a sharp distinction between measures and dimensions and does not solve the problem of being able to categorize on a measure.

- **Support for multiple hierarchies along each dimension**. For instance, Example 2.1 illustrates the type-category hierarchy on products (of interest to a consumer analyst). An alternative hierarchy is one based on which company manufactures the product and who owns that company, namely, $product \rightarrow manufacturer \rightarrow parent\ company$ (of interest to a stock market analyst). Roll-ups/drill-downs can be on either of the hierarchies.

- **Support for computing ad-hoc aggregates**. That is, aggregates other than those originally prespecified should be computable. For instance, for each product both the total sales and the average sales are interesting numbers.

- **Support for a query model in place of one-operation-at-a-time computation model**. Currently, a user operates on a cube once and obtains the resulting cube. Then the user makes the next operation. However, not all the intermediate cubes are of interest to the user. A set of basic operators that have well defined semantics enable this computation to be replaced by a query model. Thus, having tools to compose operators allows complex multidimensional queries to be built and executed faster than having the user specify each step. This approach is also more declarative and less operational.

## 2.4 Related Research

Data models developed in the context of temporal, spatial and statistical databases also incorporate dimensionality and hence have similarities with our work.

In temporal databases [TCG$^+$93], rows and columns of a relational table are viewed as two dimensions and "time" adds a third dimension forming what is called the "time cube". This cube is different from a cube in our model where dimensions correspond to arbitrary attributes and all dimensions are treated uniformly without attaching any fixed connotation with any one of them.

The modeling efforts in spatial databases [Gut94] mostly concentrate on representing arbitrary geometric objects (points, lines, polygons, regions *etc.*) in multidimensional space. By viewing OLAP data as points in the multidimensional space of attributes, one could draw analogies between the two models. But the operations central to spatial databases ("overlap", "containment", etc.) are quite different from the common OLAP operations ("roll-up", "drill-down", "joins" *etc.*). However, the multi-dimensional indexing structures developed for spatial databases (see [Gut94]) may be useful in developing efficient implementations of OLAP databases.

Statistical databases also address some of the same concerns as OLAP databases. However, models in the statistical database literature [Mic92] [CL94] have been primarily concerned with extending existing data models (mostly relational) for representing summaries and supporting operations for statistical data analysis. In contrast, our objective has been to develop a model and a set of basic operations that abstract the analysts view of enterprise data. In statistical databases, category (dimensions) and summaries (measures) are treated quite differently, whereas we have strived to treat dimensions and measures uniformly. For instance, [MERS92] describes an

S-algebra with operations like "S-union", "S-selections", "S-aggregation" that are similar to some of the operators that we define. However, since we treat measures and dimensions symmetrically, we have additional operators that are unique to our model. Irrespective of these differences, OLAP databases will benefit from implementation techniques developed in statistical databases, particularly related to aggregation views (see [Sho82] [STL89]).

Concurrent to our work, [GLS96] proposed a model for tabular data that embeds both the relational and the multidimensional data model. Their model also allows measures and dimensions to be interchanged like our model. However, their model does not state how to handle aggregations, an operation that is fundamental to OLAP databases.

## 3 Data Model

We now outline our proposed multidimensional data model and operations that capture the functionality currently provided by multidimensional database products and the additional desired functionality listed above. Our design was driven by the following key considerations:

- Treat dimensions and measures symmetrically.
- Strive for flexibility (multiple hierarchies, adhoc aggregates).
- Keep the number of operators small.
- Stay as close to relational algebra as possible. Make operators translatable to SQL.

In our logical model, data is organized in one or more multidimensional cubes. A cube has the following components:

- $k$ dimensions, and for each dimension a name $D_i$, a domain $dom_i$ from which values are taken.
- Elements defined as a mapping $E(C)$ from $dom_1 \times \ldots \times dom_k$ to either an $n$-tuple, 0, or 1. Thus, $E(C)(d_1, \ldots, d_k)$ refers to the element at "position" $d_1, \ldots, d_k$ of cube C. Note, the $d_i$s refer to values not positions per se. Therefore, our model does not require the dimensions to have a ranked, discrete domain.
- An $n$-tuple of names that describes the $n$-tuple element of the cube.

The elements of a cube can be either 0, 1, or an $n$-tuple $< X_1, \ldots, X_n >$. If the element corresponding to $E(C)(d_1, \ldots, d_k)$ is 0 then that combination of dimension values does not exist in the database. A 1 indicates the existence of that particular combination. Finally, an $n$-tuple indicates that additional information is available for that combination of dimension values. If any of the elements of a cube is a 1 then none of the elements can be a $n$-tuple and vice-versa. We represent only those values along a dimension of a cube for which at least one of

the elements of the cube is not 0. If all the elements of a cube are 0 then the cube is *empty*. Additionally, if domain $dom_i$ of dimension $D_i$ has no values then too the cube is considered to be empty.

In our model, no distinction is made between measures and dimensions. Thus, in Example 2.1, *sales* is just another dimension. Note that this is a logical model and does not force any storage mechanism. Thus, a cube in our data model may have more logical dimensions than the number of dimensions used to physically store the cube in a multidimensional storage system.

### 3.1 Operators

We now discuss our multidimensional operators. We illustrate the operators using a 2-D subset of the cube introduced in Example 2.1. We omit the *supplier* dimension and display in Figure 2 only the *product*, *date*, and *sales* dimensions. Note, *sales* is not a measure but another dimension, albeit only logical, in the model.
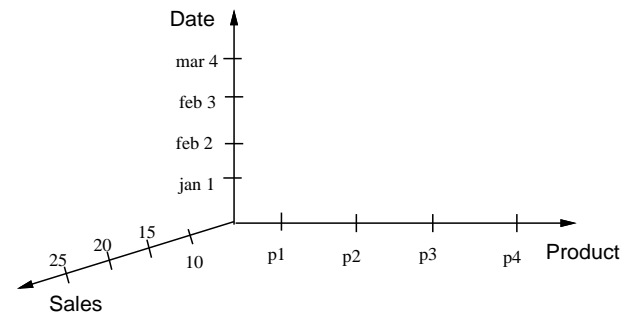


**Figure 2. Logical cube wherein *sales* is a dimension (omitting the 1/0's)**

To operate on the logical cube, the *sales* dimension may have to be folded into the cube such that sales values seem determined by the *product* and *date* dimensions. We describe later how this is achieved. For now, we will use the cube with *sales* values as the sole member of the elements of the cube. Thus, the value $< 15 >$ for "$date = mar\ 4$" and "$product = p1$" in Figure 3 indicates that in the logical cube of Figure 2 the element corresponding to "$date = mar\ 4$", "$product = p1$", and "$sales = 15$" is "1". We show the metadata description of the elements as an annotation in the cube. Thus, $<sales>$ in Figure 3 indicates that each element in the cube is a sales value.

**Notation** We define the operators using a cube $C$ with $k$ dimensions. We refer to the dimensions as $D_1, \ldots, D_k$. We use $D_i$ to refer also to the domain of dimension $D_i$ if the context makes the usage clear; otherwise we refer to the domain of dimension $D_i$ as $dom_i(C)$. We use lower
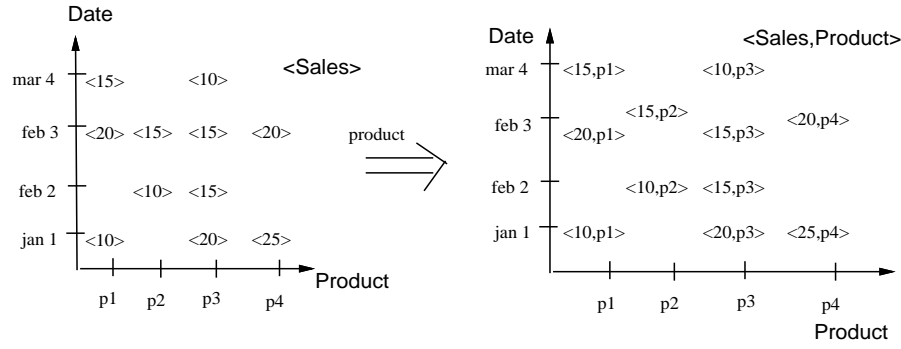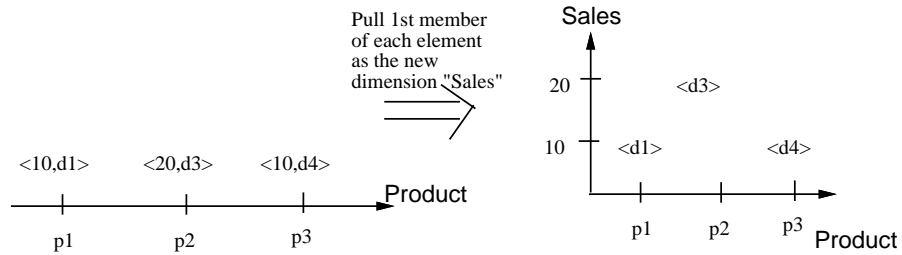
**Figure 3. The push operation on dimension** *product*



**Figure 4. Pull first member of each element as dimension** *sales*

case letters like $a$, $b$, $c$ to refer to constants.

Dimension values in our data model functionally determine elements of the cube. As a result of an application of an operation, more than one element value may be mapped to the same element (*i.e.* the same combination of values of dimension attributes) of the answer cube. These element values are combined into one value to maintain functional dependency by specifying what we call an *element combining function*, denoted as $f_{elem}$.

We also sometime merge values along a dimension. We call functions used for this purpose *dimension merging functions*, denoted as $f_{merge}$.

**Push**

The push operation (see Figure 3) is used to convert dimensions into elements that can then be manipulated using function $f_{elem}$. This operator is needed to allow dimensions and measures to be treated uniformly.

Input: $C$, $D_i$.

Output: $C$ with each non-0 element extended by an additional member, the value of dimension $D_i$ for that element.

Mathematically: $\text{push}(C, D_i) = C_a$.
$E(C_a)(d_1, \ldots, d_k)$ = $g$ $\oplus$ $d_i$ where $g$ = $E(C)(d_1, \ldots, d_k)$. The operator $\oplus$ is defined to be 0 if $g = 0$, it is $< d_i >$ if $g = 1$, and in all other cases it

concatenates $g$ and $< d_i >$.

**Pull**

This operation is the converse of the push operator. Pull creates a new dimension for a specified member of the elements. The operator is useful to convert an element into a dimension so that the element can be used for merging or joining. This operator too is needed for the symmetric treatment of dimensions and measures.

Input: $C$, new dimension name $D$, integer $i$.

Output: $C_a$ with an additional dimension $D$ that is obtained by pulling out the $i^{th}$ element of each element of the matrix.

Constraint: all non-0 elements of $C$ are $n$-tuples because each non-0 element need at least one member to enable the creation of a new dimension.

Mathematically: $\text{pull}(C, D, i) = C_a$, $1 \leq i \leq n$.
$D$ becomes the $k + 1^{st}$ dimension of the cube.
$dom_{k+1}(C_a) = \{e | e \text{ is } i^{th} \text{ member of some } E(C)(d_1, \ldots, d_k)\}$.
$E(C_a)(d_1, \ldots, d_k, e_i) =< e_1, \ldots, e_{i-1}, e_{i+1}, \ldots, e_n >$
    if $E(C)(d_1, \ldots, d_k) = < e_1, \ldots, e_i, \ldots, e_n >$,
$E(C_a)(d_1, \ldots, d_k, e_i) = 0$, otherwise.
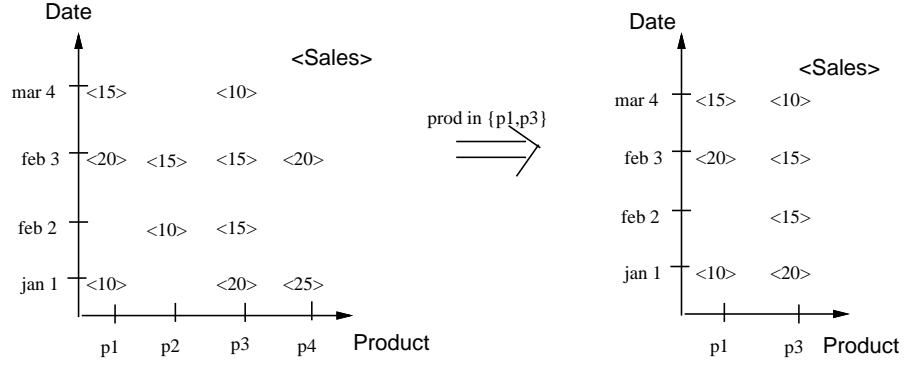Note, if $n = 1$ then elements of $C_a$ are '1's or '0's.

**Figure 5. The restriction operation**

## Destroy Dimension

Often the dimensionality of a cube needs to be reduced. This operator removes a dimension $D$ that has in its domain a single value. The presence of a single value implies that for the remaining $k - 1$ dimensions, there is a unique $k - 1$ dimensional cube. Thus, if dimension $D$ is eliminated then the resulting $k - 1$ dimensional space is occupied by this unique cube.

Input: $C$, dimension name $D_i$.

Output: $C_a$ with dimension $D_i$ absent.

Constraint: $D_i$ has only one value, say $v$

Mathematically: $destroy(C, D_i) = C_a$.
$C_a$ has $k - 1$ dimensions, $D_1 \ldots D_{i-1}, D_{i+1}, \ldots, D_k$,
$E(C_a)(d_1 \ldots d_{i-1}, d_{i+1}, \ldots, d_k) = E(C)(d_1, \ldots, d_k)$.

A dimension that has multiple values cannot be directly destroyed because then elements would no longer be functionally determined by dimension values. A multi-valued dimension is destroyed by first applying a merge operation (described later) and then applying the above operation. Note that, if $k = 1$ we will get a zero dimensional cube or a scalar as a result of the destroy operation.

## Restriction

The restrict operator operates on a dimension of a cube and removes the cube values of the dimension that do not satisfy a stated condition. Figure 5 illustrates an application of restriction. Note that this operator realizes slicing/dicing of a cube in multidimensional database terminology.

Input: Cube $C$ and predicate $P$ defined on $D_i$.

Output: New cube $C_a$ obtained by removing from $C$ those values of dimension $D_i$ that do not satisfy the predicate $P$. We have a more general notion of predicate $P$ that can be evaluated on a set of values and not on just a single value. Thus, $P$ can be either of the form "greater than 5" that is evaluated on single values at-a-time or be of the form "top 5 values" that is evaluated on the entire domain $D_i$ and outputs a set of values. If no element of dimension $D_i$ satisfies $P$ then $C_a$ is an empty cube.

Mathematically: $restrict(C, D_i, P) = C_a$.
$dom_j(C_a) = dom_j(C)$ if $1 \leq j \leq k$ & $j \neq i$
$\qquad = P(dom_j(C))$, otherwise.
$E(C_a)(d_1, \ldots, d_k) = E(C)(d_1, \ldots, d_k)$.

## Join

The join operation is used to relate information in two cubes. The result of joining a $m$-dimensional cube $C$ with an $n$-dimensional cube $C1$ on $k$ dimensions, called *joining dimensions*, is cube $C_a$ with $m + n - k$ dimensions. Each joining dimension $D_i$ of $C$ combines with exactly one dimension $D_{x_i}$ of $C1$ to get resulting dimension $D_i$ of $C_a$ as follows: For each joined dimension, two mapping functions are used for mapping the corresponding joining dimension of $C$ and $C1$ to the resulting dimension of $C_a$. The elements of $C_a$ are then formed via a function $f_{elem}$ that combines all elements of $C$ and $C1$ that get mapped to the same element of $C_a$.

Figure 6 illustrates cube $C$ joining with cube $C1$ on dimension $D_1$ (mapping function is identity). Dimension $D_1$ of the resulting cube has only two values. The function $f_{elem}$ divides the element value from cube $C$ by the element value from $C1$; if either element is 0 then the resulting element is also 0. Values of result dimension that have only 0 elements corresponding to them are eliminated from $C_a$ (like values 0 and 3 for dimension $D_1$).

Input: $C$ with dimensions $D_1 \ldots D_m$ and $C1$ with dimensions $D_{m-k+1} \ldots D_n$. $D_{m-k+1}, \ldots, D_m$ are the join dimensions (without loss of generality). $k$ mapping functions, $f_{m-k+1}, \ldots, f_m$ defined over dimensions $D_{m-k+1}, \ldots, D_m$ of $C$ and $k$ map-
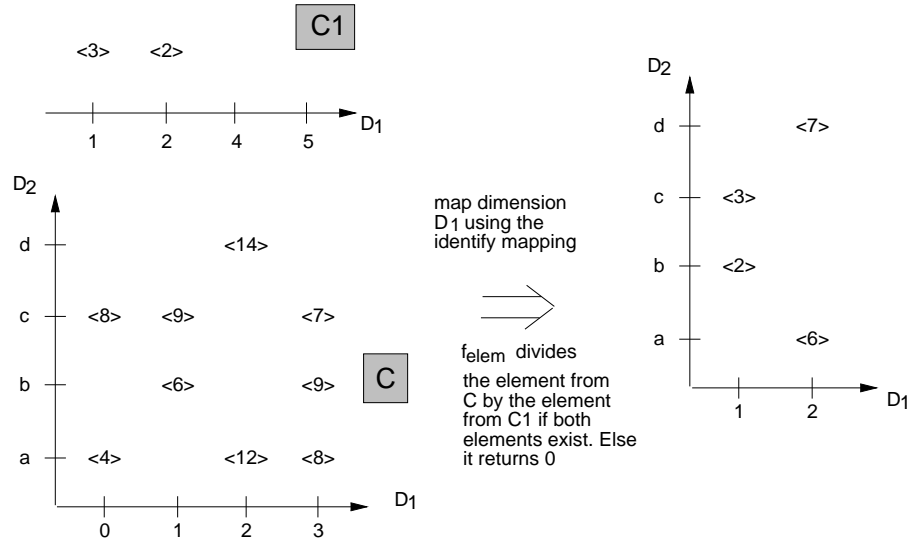
C1

<3>  <2>

1  2  4  5  D1

D2

d  <14>

c  <8>  <9>  <7>

b  <6>  <9>  C

a  <4>  <12>  <8>

0  1  2  3  D1

map dimension D1 using the identify mapping

$f_{elem}$ divides the element from C by the element from C1 if both elements exist. Else it returns 0

D2

d  <7>

c  <3>

b  <2>

a  <6>

1  2  D1

**Figure 6. Joining two cubes**

ping function $f'_{m-k+1}, \ldots, f'_m$ defined over dimensions $D_{m-k+1}, \ldots, D_m$ of $C1$. Mapping $f_i$ applied to value $v \in dom_i(C)$ produces values for dimension $D_i$ in $C_a$. Similarly $f'_i$ applied to $v' \in dom_i(C1)$ produces values for dimension $D_i$ in $C_a$. Also needed is a function $f_{elem}$ that combines sets of elements from $C$ and $C1$ to output elements of $C_a$.

Output: $C_a$ with $n$ dimensions, $D_1 \ldots D_n$. Multiple elements in $C$ and $C1$ could get mapped to the same element in $C_a$. All elements of $C$ and $C1$ that get mapped to the same point in $C_a$ are combined by function $f_{elem}$ to produce the output element of $C_a$. If for some value $v$ of dimension $D_i$, the elements $E(C_a)(x_1, \ldots, v, x_{i+1}, \ldots, x_n)$ is 0 for all values of the other dimensions, then $v$ is not included in dimension $D_i$ of $C_a$.

Mathematically:
$$\text{join}(C, C1, [f_{m-k+1}, \ldots, f_m, f'_{m-k+1}, \ldots, f'_m], f_{elem})$$
$$= C_a$$
$$dom_i(C_a) = dom_i(C) \text{ if } 1 \le i \le m-k.$$
$$= dom_i(C1) \text{ if } m+1 \le i \le n.$$
$$= \{d^a | d^a \in f_i(d), d \in dom_i(C) \text{ OR}$$
$$d^a \in f'_i(d'), d' \in dom_i(C1)\}.$$
$$E(C_a)(d_1, \ldots, d^a_{m-k}, \ldots, d^a_m, \ldots, d_n) = f_{elem}(\{t1\}, \{t2\})$$
such that
$$t1 = E(C)(d_1, \ldots, d_{m-k}, \ldots, d_m),$$
$$t2 = E(C1)(d'_{m-k}, \ldots, d'_m, d_{m+1}, \ldots, d_n),$$
$$d^a_i \in f_i(d_i) \text{ OR } d^a_i \in f'_i(d'_i)$$

We defined above a general notion of the join operator that covers several important special cases. Notable amongst these are: **cartesian product**, **natural join**, **union**, **merge**

and **associate**. In the case of **cartesian product**, the two cubes have no common joining dimension. In the case of **natural join** the mapping function is identity and the $f_{elem}$ function returns "0" whenever one of the elements is '0'. The sub-cases union and merge are discussed later.

**Associate** is an especially useful sub case in OLAP applications for computations like "express each month's sale as a percentage of the quarterly sale." Associate is asymmetric and requires that each dimension of $C1$ be joined with some dimension of $C$. Figure 7 illustrates associating cube $C1$ with $C$ where *month* dimension of $C1$ and *date* dimension of $C$ are joined by mapping them to the *date* dimension of $C_a$. Similarly, *category* and *product* are joined by mapping them to *product* of $C_a$. For dimension *month*, each month is mapped to all the dates in that month. For dimension *category*, value $cat1$ is mapped to *products* $p1$ and $p2$, and $cat2$ is mapped to $p3$ and $p4$. For dimensions *date* and *product* the identity mapping is used. Function $f_{elem}$ divides the element value from cube $C$ by the element value from $C1$; if either element is 0 then the resulting element is also 0. Note, value *mar4* is eliminated from $C_a$ because all its corresponding elements are 0.

**Merge**

The merge operation is an aggregation operation. We illustrate it in Figure 8. The figure shows how hierarchies in a multidimensional database are implemented using the merge operator. Intuitively, a dimension merging function is used to map multiple product names into one or more categories and another function is used to map individual dates into their corresponding month. Thus, multiple elements on each dimension are merged to produce a di-
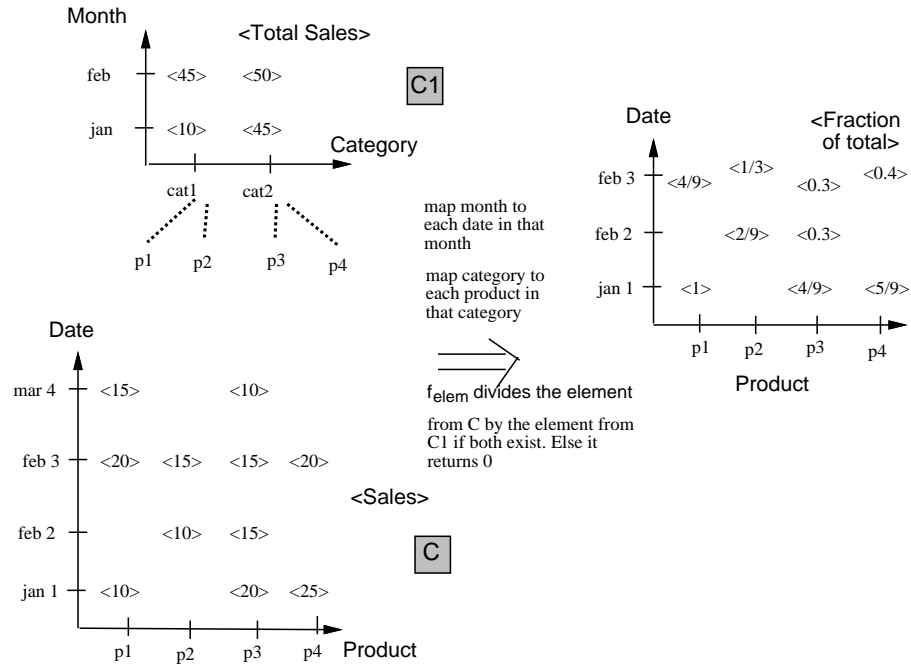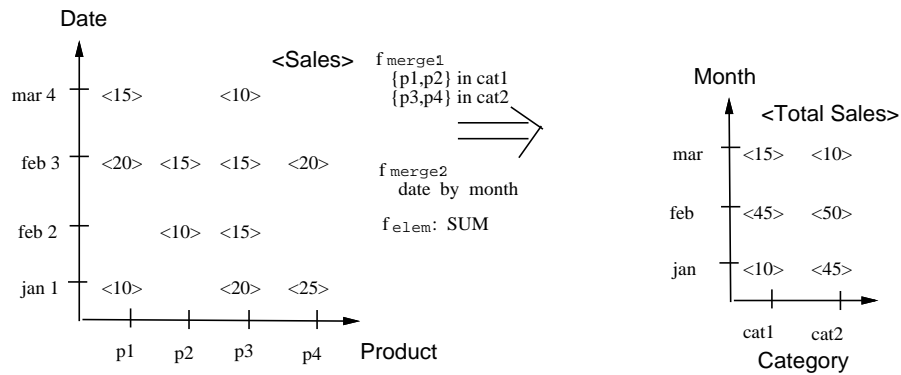
**Figure 7. Associating two cubes**



**Figure 8. Merging dimensions** *date* **and** *product* **using** $f_{elem} = sum$

mension with a possibly smaller domain. As a result of merging each dimension, multiple elements in the original cube get mapped to the same element in the new cube. An element combining function is used to specify how these multiple elements are aggregated into one. In the example in Figure 8, the aggregation function $f_{elem}$ is **sum**.

In general, the dimension merging function might be a one-to-many mapping that takes an element in the lower level into *multiple* elements in the higher level of hierarchies. For instance, a $1 \rightarrow n$ mapping can be used to merge a product belonging to $n$ categories to handle multiple hierarchies.

Input: $C$, function $f_{elem}$ for merging elements and $m$ (dimension, function) pairs. Without loss of generality, as-

sume that the $m$ pairs are $[D_1, f_1], \ldots, [D_m, f_m]$

Output: Cube $C_a$ of same dimensionality as $C$. Dimension $D_i$ is merged as per function $f_i$. An element corresponding to the merged elements is aggregated as per $f_{elem}$.

Mathematically:
merge$(C, \{[D_1, f_1], \ldots, [D_m, f_m]\}, f_{elem})$=$C_a$.
$dom_i(C_a) = \{f_i(e) | e \in dom_i(C)\}$ if $1 \le i \le m$,
$\qquad = dom_i(C)$, otherwise.
$E(C_a)(d_1, \ldots, d_k) = f_{elem}(\{t | t = E(C)(d'_1, \ldots, d'_k)$
$\qquad$ where for $f_i(d'_i) = d_i$ if $1 \le i \le m$ else $d'_i = d_i\})$.

A special case of the merge operator is when all the

merging functions are identity. In this case, the merge operator can be used to *apply* a function $f_{elem}$ to each element of a cube.

**Remark** The merge operator is strictly not part of our basic set of operators. It can be expressed as a special case of the self-join of a cube using $f_{merge}$ transformation functions on dimensions being merged and identity transformation functions for other dimensions. We chose to separately define merge because it is a unary operator unlike the binary join and also for performance reasons.

## 4 Discussion

The reader may have noticed similarities in the operators proposed and relational algebra [Cod70]. It is by design. One of our goals was to explore how much of the functionality of current multidimensional products can be abstracted in terms of relational algebra. By developing operators that can be translated into SQL (see [AGS96] for translations), our hope was to create a fast path for providing OLAP capability on top of relational database systems. We must hasten to add that we are not arguing against specialized OLAP engines—we believe the design and prototyping of such engines is a fruitful research direction. We are also not suggesting that simply translating these operators into SQL would be sufficient for providing OLAP capabilities in relational database systems. However, it does point to directions in which optimization techniques and storage structures in relational database systems need to evolve.

The goal of treating dimensions and measures symmetrically had a permeating influence in our design. It is a functionality either not present or poorly supported in current multidimensional database products. Its absence causes expensive schema redesign when an unanticipated need arises for aggregation along an attribute that was initially thought to be a measure. In hindsight, the push and pull operations may appear trivial. However, their introduction was the key that made the symmetric treatment of dimensions and measures possible.

The reader may argue with the way we have chosen to incorporate order-based information into our algebra. We rely on functions for this purpose, which implies that the system may not be able to use this information in optimizing queries. We debated about allowing a native order to be specified with each dimension and providing ordering operators. We decided against it because of the large number of such operators and because the semantics gets quite complex when there are multiple hierarchies along a dimension. In a practical implementation of our model, it will be worthwhile to allow a default order to be specified with each dimension and make system aware of some built-in ordering functions such as "first $n$". The same holds for providing the knowledge of some built-in aggregate functions.

The reader may also notice the absence of direct analogs of relational projection, union, intersection, and difference. These operations can be expressed in terms of our proposed operators as follows:

**Projection** The projection of a cube is computed by merging each dimension not included in the projection and then destroying the dimension. A $f_{elem}$ function specifying how multiple elements are combined is needed as part of the specification of the projection.

**Union** Two cubes are union-compatible if (i) they have the same number of dimensions; and (ii) for each dimension $D_i$ in $C$, dimension $D_i$ in $C1$ is of the same domain type. Union is computed by joining the two cubes using the identity transformation functions for each dimension of each cube and by choosing a function $f_{elem}$ that produces a non-empty element for element $e$ in $C_a$ whenever an element from either of the two cubes is mapped into $e$. Dimension $D_i$ in the resulting cube has as its values the union of the values in $dom_i(C)$ and in $dom_i(C1)$.

**Intersect** The intersection of two union-compatible cubes is computed by joining the cubes through the identity mapping that effectively retains only those dimension values that are present in both cubes. Thus, function $f_{elem}$ makes non-0 an element for point $p$ in $C_a$ only if elements from both cubes are mapped into $p$.

**Difference** The difference of two union-compatible cubes $C1$ and $C2$ is expressed as an intersection of $C1$ and $C2$ followed by a union of the result with $C1$. The $f_{elem}$ function for combining two elements for the intersection steps discards the value of the element for $C1$ and retains $C2$'s element. The $f_{elem}$ function for combining two elements for the union step saves the value of $C1$'s element if the two elements are different and makes the result 0 if they are identical[1].

### 4.1 Expressive Power

Our algebra can be seen to be at least as powerful as relational algebra [Cod70]. A precise circumscription of the expressive power of the proposed data model is an open problem. A related interesting open question pertains to defining a formal notion of completeness for multidimensional database queries and evaluating how complete our algebra is with respect to that metric. We take an empirical approach and discuss below how the current high-level multidimensional operations can be built using our proposed operators.

---

[1]This implementation corresponds to the following semantics for $C1 - C2$: $E(C_a)(d_1, \ldots, d_k)$ equals 0 if $E(C2)(d_1, \ldots, d_k) = E(C1)(d_1, \ldots, d_k)$; it is $E(C1)(d_1, \ldots, d_k)$ otherwise. Another alternative semantics could be that $E(C_a)(d_1, \ldots, d_k)$ equals 0 if $E(C2)(d_1, \ldots, d_k) \neq 0$, and $E(C1)(d_1, \ldots, d_k)$ otherwise. This semantics can be implemented by a small change in the $f_{elem}$ function used in the union step.

**Roll-up**   Roll-up is a merge operation that needs one dimension merging function and one element combining function. If a hierarchy is specified on a dimension then the dimension merging function is defined implicitly by the hierarchy. The elements corresponding to merged values on the dimension are combined using the user-specified element combining function like *SUM*.

**Drill-down**   This operator is actually a binary operation even though most current multidimensional database products make it seem like a unary operation. Consider computing the sum $X$ of 10 values. Drill-down from $X$ to the underlying 10 values is possible in infinite ways. Thus, the underlying 10 values have to be known. That is, the aggregate cube has to be joined (actually associated) with the cube that has detailed information. Continuing with our analogy, to drill down from $X$ to its constituents the database has to keep track of how $X$ was obtained and then associate $X$ with these values. Thus, if users merge cubes along stored paths and there are unique path down the merging tree, then drill down is uniquely specified. By storing hierarchy information and by restricting single element merging functions to be used along each hierarchy, drill-down can be provided as a high-level operation on top of associate.

**Star Join**   In a star join [Eri95], a large detail "mother" table $M$ is joined with several small "daughter" tables that describe join keys in the mother table. A star join denormalizes the mother table by joining it with its daughter tables after applying selection conditions to their descriptive attributes. We describe how our operators capture a star join when $M$ has one daughter table $F_1$ that describes the join key field $D$ of $M$. Table $F_1$ can be viewed as a one-dimensional cube, $C_1$ with the join key field $D$ as the dimension and all the description fields pulled in as elements. A restriction on a description attribute $A$ of table $F_1$ corresponds to a function application to the elements of $C_1$. Restrictions on the join key attribute translate to restrictions on dimension $D$ of $C_1$. The join between $M$ and $F_1$ is achieved by associating the mother cube with the daughter cube on the key dimension $D$ using the identity mapping function. The description of each key value is pulled in from the daughter cube into the mother cube via the $f_{elem}$ function.

**Expressing a dimension as a function of other dimensions**   This functionality is basic in spread sheets. We can create a new dimension $D$ expressed as a function, $f$ of another dimension $D'$ by first pushing $D'$ into the cube elements, then modifying the cube elements by applying function $f$ and finally pulling out the corresponding member of the cube element as a new dimension $D$.

## 4.2   Example queries

This section illustrates how to express some of the queries of Example 2.2 using our operators. Assume we have a cube $C$ with dimensions product, month, supplier and element sales.

*For supplier "Ace" and for each product, give the fractional increase in the sales in January 1995 relative to the sales in January 1994.*

**Restrict** supplier to "Ace" and dates to "January 1994 or January 1995". **Merge** date dimension using an $f_{elem}$ that combines sales as $(B - A)/A$ where $A$ is the sale in Jan 1994 and $B$ is the sale in Jan 1995.

*For each product give its market share in its category this month minus its market share in its category in October 1994.*

**Restrict** date to "October 1994 or current month". **Merge** supplier to a single point using sum of sales as the $f_{elem}$ function to get $C1$. **Merge** product dimension to category using sum as the $f_{elem}$ function to get in $C2$ the total sale for the two months of interest. **Associate** $C1$ and $C2$, mapping a category in $C2$ to each of its products in $C1$. The identity mapping is used for the Month dimension. Function $f_{elem}$ divides the element from $C1$ by the element from $C2$ to get the market share. For the resulting cube, **Merge** dimension month to a single point using a $f_{elem}$ function $(A - B)$ where $A$ is the market share for "this" month and $B$ is the market share in October 1994.

*For each product category, select total sales this month of the product that had highest sales in that category last month.*

**Restrict** dimension month to "last" month. **Merge** supplier to a single point using sum of sales as the $f_{elem}$ function. **Push** product dimension resulting in 2-tuple elements with <Sale and product>. **Merge** product to category using $f_{elem}$ function that retains an element if it has the "maximum" sales. **Pull** product into the category dimension (over-riding the category dimension, this can be easily done using our basic operators). Let the resulting cube be $C1$. This cube has the highest sales value for each element for "last" month. **Restrict** $C$ on dimension date to "this" month, **Merge** supplier to a single point using sum of sales as the $f_{elem}$ function and **associate** it with $C1$ on the product dimension using $f_{elem}$ function that only outputs the element of $C$ when it is the same as the corresponding elements from $C1$ (otherwise returns 0).

*Select suppliers for which the total sale of every product increased in each of last 5 years.*

**Restrict** to months of last 6 years. **Merge** month to year. **Merge** years to a single point using a $f_{elem}$ function that maps the six sales values to "1" if sales values are increasing, "0" otherwise. **Merge** product to a point where $f_{elem}$ function is "1" if and only if all its arguments are "1".

# 5 Conclusions and Future Work

This paper introduced a data model and a set of algebraic operations that unify and extend the functionality provided by current multidimensional database products. As illustrated in Section 4.1, the proposed operators can be composed to build OLAP operators like roll-up, drill-down, star-join and many others. In addition, the model provides symmetric treatment to dimensions and measures. The model also provides support for multiple hierarchies along each dimension and support for adhoc aggregates. Absence of these features in current products results in expensive schema redesign when an unanticipated need arises for a new aggregation or aggregation along an attribute that was initially thought to be a measure.

The proposed operators have several desirable properties. They have well-defined semantics. They are minimal in the sense that none can be expressed in terms of others nor can any one be dropped without sacrificing functionality. Every operator is defined on cubes and produces as output a cube. That is, the operators are closed and can be freely composed and reordered. This allows the inefficient one-operation-at-a-time approach currently in vogue to be replaced by a query model and makes multidimensional queries amenable to optimization.

The proposed operators enable the logical separation of the frontend user interface from the backend that stores and executes queries. They thus provide an algebraic API that allows the interchange of frontends and backends. The operators are designed to be translated into SQL. Thus, they can be implemented on either a relational system or a specialized engine.

For future, on the modeling side, work is needed to incorporate duplicates and NULL values in our model. We believe that the duplicates can be handled by treating elements of the cube as pairs consisting of an arity and a tuple of values. The arity gives the number of occurrences of the corresponding combination of dimensional values. NULLs can be represented by allowing for a NULL value for each dimension. Details of these extensions and other possible alternatives require further investigation.

On the implementation side, there are interesting research problems for implementing our model on top of a relational system as well as within a specialized engine. Although each of the proposed operators can be translated into a SQL query, simply executing this translated SQL on a relational engine is likely to be quite inefficient. Corresponding to a multidimensional query composed of several of these operators, we will get a sequence of SQL queries that offers opportunity for multi-query optimization. It needs to be investigated whether the known techniques (e.g. [SG90]) will suffice or do we need to develop new techniques. Similarly, there is opportunity for research in storage and access structures and materialized views.

# References

[AAD+96] S. Agarwal, R. Agrawal, P.M. Deshpande, A. Gupta, J.F. Naughton, R. Ramakrishnan, and S. Sarawagi. On the computation of multidimensional aggregates. In *Proc. of the 22nd Int'l Conference on Very Large Databases*, pages 506–521, Mumbai (Bombay), India, September 1996.

[AGS96] Rakesh Agrawal, Ashish Gupta, and Sunita Sarawagi. Modeling multidimensional databases. Research Report, IBM Almaden Research Center, San Jose, California, 1996. Available from `http://www.almaden.ibm.com/cs/quest`.

[Arb] Arbor Software Corporation, Sunnyvale, CA. *Multidimensional Analysis: Converting Corporate Data into Strategic Information.* `http://www.arborsoft.com`.

[CCS93] E. F. Codd, S. B. Codd, and C. T. Salley. Beyond decision support. *Computerworld*, 27(30), July 1993.

[CL94] R. Cicchetti and L. Lakhal. Matrix relation for statistical database management. In *Proc. of the Fourth Int'l Conference on Extending Database Technology (EDBT)*, March 1994.

[Cod70] E.F. Codd. A relational model for large shared data banks. *Comm. ACM*, 13(6):377–387, 1970.

[Cod93] E. F. Codd. Providing OLAP (on-line analytical processing) to user-analysts: An IT mandate. Technical report, E. F. Codd and Associates, 1993.

[Col95] George Colliat. OLAP, relational, and multidimensional database systems. Technical report, Arbor Software Corporation, Sunnyvale, CA, 1995.

[Eri95] Christopher G. Erickson. Multidimensionalism and the data warehouse. In *The Data Warehousing Conference*, Orlando, Florida, February 1995.

[Fin95] Richard Finkelstein. MDD: Database reaches the next dimension. *Database Programming and Design*, pages 27–38, April 1995.

[GBLP96] J. Gray, A. Bosworth, A. Layman, and H. Pirahesh. Data cube: A relational aggregation operator generalizing group-by, cross-

tabs and sub-totals. In *Proc. of the 12th Int'l Conference on Data Engineering*, pages 152–159, 1996.

[GHRU97] Himanshu Gupta, Venky Harinarayan, Anand Rajaraman, and Jeffrey D. Ullman. Index selection for OLAP. In *Proc. of the 13th Int'l Conference on Data Engineering*, Birmingham, U.K., April 1997.

[GLS96] M. Gyssens, L.V.S. Lakshmanan, and I.N. Subramanian. Tables as a paradigm for querying and restructuring. In *Proceedings of the ACM Symposium on Principles of Database Systems (PODS)*, 1996.

[Gut94] R. H. Guting. An introduction to spatial database systems. *VLDB Journal*, 3(4):357–399, 1994.

[HRU96] V. Harinarayan, A. Rajaraman, and J.D. Ullman. Implementing data cubes efficiently. In *Proc. of the ACM SIGMOD Conference on Management of Data*, June 1996.

[IRI] IRI Software, Information Resources Inc., Waltham, MA. *OLAP: Turning Corporate Data into Business Intelligence*. http://www.infores.com.

[JS96] T. Johnson and D. Shasha. Hierarchically split cube forests for decision support: description and tuned design, 1996. Working Paper.

[KS94] R. Kimball and K. Strehlo. What's wrong with SQL. *Datamation*, June 1994.

[MERS92] L. Meo-Evoli, F.L. Ricci, and A. Shoshani. On the semantic completeness of macro-data operators for statistical aggregations. In *Proceedings of the Sixth International Working Conference on Scientific and Statistical Database Management*, 1992.

[Mic] Microstrategy Inc., Vienna, VA 22182. *True Relational OLAP*. http://www.microstrategy.com.

[Mic92] Z. Michalewicz. *Statistical and Scientific Databases*. Ellis Horwood, 1992.

[OLA96] The OLAP Council. *MD-API the OLAP Application Program Interface Version 0.5 Specification*, September 1996.

[Rad95] Neil Raden. Data, data everywhere. *Information Week*, pages 60–65, October 30 1995.

[SDNR96] A. Shukla, P.M. Deshpande, J.F. Naughton, and K. Ramasamy. Storage estimation for multidimensional aggregates in the presence of hierarchies. In *Proc. of the 22nd Int'l Conference on Very Large Databases*, pages 522–531, Mumbai (Bombay), India, September 1996.

[SG90] T. Sellis and S. Ghosh. On the multiple-query optimization problem. *IEEE Transactions on Knowledge and Data Engineering*, 2(2):262–266, 1990.

[Sho82] A. Shoshani. Statistical databases: Characteristics, problems and some solutions. In *Proceedings of the Eighth International Conference on Very Large Databases (VLDB)*, pages 208–213, Mexico City, Mexico, September 1982.

[SR96] B. Salzberg and A. Reuter. Indexing for aggregation, 1996. Working Paper.

[STL89] J. Srivastava, J.S.E. Tan, and V.Y. Lum. TBSAM: An access method for efficient processing of statistical queries. *IEEE Transactions on Knowledge and Data Engineering*, 1(4), 1989.

[TCG+93] A.U. Tansel, J. Clifford, S. Gadia, S. Jajodia, A. Segev, and R. Snodgrass. *Temporal Databases: Theory, Design, and Implementation*. Benjamin/Cummings, 1993.