

# Querying Multidimensional Databases\*

Luca Cabibbo and Riccardo Torlone

Dipartimento di Informatica e Automazione

Università di Roma Tre

Via della Vasca Navale, 79 — I-00146 Roma, Italy

E-mail: {cabibbo,torlone}@inf.uniroma3.it

**Abstract.** Multidimensional databases are large collections of data, often historical, used for sophisticated analysis oriented to decision making. This activity is supported by an emerging category of software technology, called On-Line Analytical Processing (OLAP). In spite of a lot of commercial tools already available, a fundamental study for OLAP systems is still lacking. In this paper we introduce a model and a query language to establish a theoretical basis for multi-dimensional data. The model is based on the notions of dimension and f-table. Dimensions are linguistic categories corresponding to different ways of looking at the information. F-tables are the constructs used to represent factual data, and are the logical counterpart of multi-dimensional arrays, the way in which current analytical tools store data. The query language is a calculus for f-tables, and as such it offers a high-level support to multi-dimensional data analysis. Scalar and aggregate functions can be embedded in calculus expressions in a natural way. We discuss on conceptual problems related with the design of multidimensional query languages, and compare our model and language with other approaches.

## 1 Introduction

The integration of database management systems with on-line analytical processing (OLAP) technology is a challenging goal of recent years [8]. In fact, while the former provide solid and efficient tools for on-line *transaction* processing, OLAP systems can support knowledge workers and decision makers in the sophisticated analysis of enterprise data. The effectiveness of this analysis is related to the ability to describe and manipulate data according to different and often independent perspectives or “dimensions.” For instance, single sales of items provide much more information to business analysis when organized into, e.g., number of items sold by category of product, geographical location, and time. Thus, we can say that OLAP technology complements database technology, in providing a multi-dimensional view of raw data and suitable tools for its analysis. Generally, these tools enable the users to: (i) define analytical equations across multiple data dimensions, possibly involving complex calculations, to represent numerous, speculative enterprise model scenarios; (ii) summarize data sets, aggregating

---

\* This work was partially supported by *Consiglio Nazionale delle Ricerche* and by *MURST*.

and disaggregating over the various dimensions; and (iii) evaluate and view the outcomes of the analysis. To understand the effect of changes in environmental factors, this process is often iterated by changing equations and parameters.

Current technology provides both OLAP data servers and front-end analysis tools. The former can be either relational systems (ROLAP) or proprietary multi-dimensional database systems (MOLAP) [9]. The latter offer interactive graphical user interfaces, usually similar to spreadsheets. While this allows the user to easily summarize and view data, spreadsheet-like environments suffer from several limitations in constructing and maintaining analytical models over the enterprise data. The main point is that these models rely on a logic that is often left implicit, leading to several problems, including redundancy and inconsistency [15]. Moreover, the integration with database technology is based on ad-hoc techniques, rather than any systematic approach. As others [12], we believe that the problem is the lack of a formal theoretical foundation.

In this paper, we propose the MultiDimensional data model and query language, as a new basis for OLAP systems. The model allows to describe the logical structure of the enterprise data according to multiple perspectives, by providing an explicit notion of *dimension*. Dimensions are the linguistic categories used to characterize the structure of data, according to a conceptual business perspective. They are organized into hierarchies of levels, corresponding to possible granularities of data. Factual data are then represented by *f-tables*, the logical counterpart of “multi-dimensional arrays” (the way in which OLAP systems store values). Values in f-tables are accessed through symbolic coordinates. The query language enables the user to express cross-dimensional analytical equations, based on logical expressions over f-tables, in a simple and declarative way. Queries make use of interpreted functions, but the language is parametric with respect to the ones chosen. A distinctive feature of our model is the use of *roll-up functions*, which describe how data are related within hierarchies of levels. Roll-up functions provide the query language with a simple and powerful mechanism to join data at different levels of aggregation.

The main contributions of the paper are the following. The presentation of a multi-dimensional database model, which provides a first step towards a logical foundation of OLAP systems. The development of a calculus-like query language, which offers a high-level support to multi-dimensional data analysis. The study of the expressiveness of the model and the query language, based on a comparison with other related approaches in the literature. In particular, we show that our model subsumes the relational data model. We then prove that our query language, with a suitable collection of functions, expresses the relational algebra, eventually extended with aggregate functions [16]. We also show that, with a limited set of interpreted functions, the query language is able to express a broad class of queries.

*Related work.* The term OLAP has been recently introduced by Codd et al. [8] to characterize the category of analytical processing over large, historical databases (data warehouses) oriented to decision making. Further discussion on OLAP, multi-dimensional analysis, and data warehousing can be found in [6,14,22,24].

A comparison between OLAP concepts and the area of statistical databases is given in [20].

An important OLAP operation is summarization of data over one or more dimensions. Klug [16] provided a first theoretical basis in this respect, by extending the relational algebra and calculus with aggregate functions, that is, interpreted functions taking a set of tuples as argument and producing a single value as result. Our approach is more general than Klug's one, since the MultiDimensional model subsumes the relational one. Furthermore we consider, in addition to aggregate functions, also scalar functions.

Many authors [5,11,19,23] claim that SQL is unsuited to data-analysis applications, since some aggregate and grouping queries are difficult to express and optimize. They thus consider the problem of extending SQL with aggregation and analysis-oriented operators that are more powerful, but also specific to particular application domains. Gray et al. [11] propose *cube* as an operator generalizing *group by*. Chatziantoniou and Ross [5] extend both SQL and the relational algebra with an operator, dealing with "aggregation variables", to succinctly express common queries, providing also a basis for improved query optimization. Rao et al. [19] consider the issue of supporting quantified queries, a class of queries that is difficult to deal with in SQL; they introduce a number of "generalized quantifiers", that is, aggregate set-predicates such as *some*, *all*, and *at-least*. Many of the features considered in such proposals can be easily expressed in our language using a limited collections of scalar and aggregate interpreted functions.

Agrawal et al. [4] have proposed a simple hypercube-based data model, and a few algebraic operators for it. This framework shares a number of characteristics and goals with ours. However, the approach is rather pragmatic, mainly oriented towards a direct SQL implementation into a relational database. Conversely, we have followed a more systematic approach. Moreover, our notion of roll-up function let us capture and describe hierarchies of levels within dimensions in a neater way, and allows us to express more easily complex aggregations in queries.

Libkin et al. [17] have defined a language for querying data organized in multi-dimensional arrays, to support the scientific computing community with database technology. The MultiDimensional model is at a different, and perhaps higher, abstraction level; our notion of f-table is indeed a "logical" counterpart of a "physical" multi-dimensional array. It should also be noted that our approach is motivated by a business context.

Gyssens et al. [12] have proposed the tabular database model, together with a complete algebraic language for querying and restructuring, as a first theoretical foundation for OLAP systems. A main difference with respect to their approach is that we introduce an explicit logical notion of dimension, allowing for multi-dimensional structures, whereas their tables are bidimensional. Their query language covers only the aspect of restructuring, whereas we allow complex computations based on formulas and functions.

**Organization.** The paper is organized as follows. The MultiDimensional data model is presented in Section 2. The associated query language is introduced

informally, by means of examples, in Section 3, and described formally in Section 4. Section 5 presents expressiveness results. Finally, Section 6 discusses further research topics.

## 2 The MultiDimensional Data Model

This section introduces the MultiDimensional data model (MD for short). The model is based on the notion of *dimension* that allows to specify multiple “ways” to look at information, according to natural business perspectives under which its analysis can be performed. Each dimension is organized in a hierarchy of *levels*, corresponding to data domains at different granularities. A MultiDimensional *scheme* consists of a set of *f-tables* that are defined with respect to particular combinations of levels. A MultiDimensional *instance* associates *measures*, which correspond to data being tracked, with *symbolic coordinates* over *f-tables*.<sup>1</sup> Finally, within a dimension, values of a finer granularity can *roll up* to (that is, can be grouped into) values of a coarser one.

*Example 1.* A marketing analyst of a chain of toy stores may organize its business data along dimensions like time, product, and location. The time dimension may be organized in levels *day*, *quarter*, *week*, and *year*, and Feb 19, 97 is an element of the *day* level. The elements of this level roll up to elements of levels *week* and *quarter*. Similarly, both *weeks* and *quarters* roll up to *years*. Note however that *weeks* do not roll up to *months*, since months do not divide evenly into weeks. In this framework, a measure can be the number of items sold by the chain. This measure could be represented by means of an *f-table* *SALES*, having symbolic coordinates on the levels *day*, *item*, and *store*: an instance of *SALES* might state the fact that on Feb 19, 97 the store Colosseum has sold 11 pieces of Lego.

### 2.1 MultiDimensional Schemes

Let us fix two disjoint countable sets of *names* and *values*. We denote by  $\mathcal{L}$  a set of names called *levels* such that: (i) each level  $l \in \mathcal{L}$  is associated with a countable set of values  $\text{DOM}(l)$ , called the *domain* of  $l$ ; and (ii) the various domains associated with different levels are pairwise disjoint.

**Definition 2 (Dimension).** A *dimension*  $d$  is a triple  $(L, \preceq, \text{R-UP})$ , where:

- $L \subseteq \mathcal{L}$  is a finite set of levels;
- $\preceq$  is a partial order defined among the levels of  $d$ . Whenever  $l_1 \preceq l_2$  we say that  $l_1$  *rolls up* to  $l_2$  or that  $l_2$  *drills down* to  $l_1$ ;

<sup>1</sup> Actually, the ‘f’ in the term ‘f-table’ has a double meaning. On one hand, it stands for ‘function’, because each f-table is indeed a function, from coordinates to measures. On the other hand, it stands also for ‘fact’, since f-tables represent a form of information that practitioners implement by means of the so-called ‘fact tables’.

- R-UP is a family of functions, called *roll-up functions*, satisfying the following conditions:
  - for each pair of levels  $l_1, l_2$  such that  $l_1 \preceq l_2$ , the roll-up function  $\text{R-UP}_{l_1}^{l_2}$  maps each element of  $\text{DOM}(l_1)$  to an element of  $\text{DOM}(l_2)$ . Whenever  $\text{R-UP}_{l_1}^{l_2}(o_1) = o_2$  we say that  $o_1$  *rolls up to*  $o_2$ , or that  $o_2$  *drills down to*  $o_1$ ;
  - given levels  $l_1, l'$ , and  $l_2$  such that  $l_1 \preceq l' \preceq l_2$ , (and thus,  $l_1 \preceq l_2$ ) the function  $\text{R-UP}_{l_1}^{l_2}$  equals the composition  $\text{R-UP}_{l'}^{l_2} \circ \text{R-UP}_{l_1}^{l'}$ . This implies that: (i) for each level  $l$ , the function  $\text{R-UP}_l^l$  is the identity on  $\text{DOM}(l)$ ; and (ii) whenever a level  $l_1$  rolls up to  $l_2$  in different ways (e.g., rolling up through either  $l'$  or  $l''$ ) then the elements of  $l_1$  roll up to elements of  $l_2$  in a consistent way.

*Example 3.* Consider Example 1. The relevant information is organized along dimensions **time**, **product**, and **location**, and involves numeric data describing sales and prices. The dimension hierarchies are depicted on top of Figure 1; note that each dimension takes the name from one of its levels (often the least upper bound of its lattice). The figure shows that, e.g., level **item** rolls up to both **category** and **brand**; because of reflexivity, **item** rolls up also to itself and, because of transitivity, it rolls up to **product**. The domain associated with the level **day** contains, among others, values Jan 5, 97, Feb 19, 97, and Mar 10, 97, all of which roll up to the element 1Q-97 of the level **quarter**. The level **store** contains values Colosseum and Navona, both of them rolling up to Rome (in level **city**) and Italy (in level **area**). The level **numeric** is a built-in level type, having as domain the rational numbers.

**Definition 4 (Scheme).** A *MultiDimensional scheme* is a pair  $(D, F)$ , where:

- $D$  is a finite set of *dimensions*;
- $F$  is a finite set of *f-table schemes* of the form  $f[A_1 : l_1, \dots, A_n : l_n] : l_0$ , where  $f$  is a name (with the condition that different f-table schemes have distinct names), each  $A_i$ , for  $1 \leq i \leq n$ , is a distinct name (called an *attribute of f*), and each  $l_i$ , for  $0 \leq i \leq n$ , is a level of some dimension in  $D$ .

*Example 5.* Figure 1 shows the MD scheme **Toys**, having two f-tables, named **SALES** and **PRICE-LIST**. Intuitively, the f-table **SALES** represents summary data for the sales of the chain in terms of pieces sold (dimension **numeric**), organized along dimensions **time** (at **day** level), **location** (at **store** level), and **product** (at **item** level). F-table **PRICE-LIST** is instead used to price the various items, assuming that prices may vary from month-to-month, and that different stores sell each item at the same price.

## 2.2 MultiDimensional Instances

**Definition 6 (Coordinate).** Let  $\mathcal{S} = (D, F)$  be a MultiDimensional scheme and  $f[A_1 : l_1, \dots, A_n : l_n] : l_0$  be an f-table scheme in  $F$ . A (*symbolic*) *coordinate*

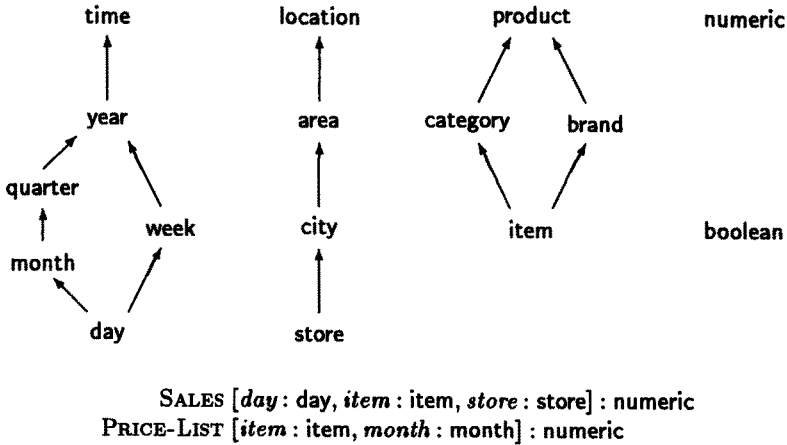


Fig. 1. The sample Toys scheme.

$\gamma$  over  $f$  is a function mapping each attribute name  $A_i$  (with  $1 \leq i \leq n$ ) to an element  $o_i \in \text{DOM}(l_i)$ . If  $\gamma$  is a coordinate over  $f$  such that  $\gamma(A_i) = o_i$ , for  $1 \leq i \leq n$ , we denote  $\gamma$  by  $[A_1 : o_1, \dots, A_n : o_n]$ .

**Definition 7 (Instance).** Let  $S = (D, F)$  be a MultiDimensional scheme and  $f[A_1 : l_1, \dots, A_n : l_n] : l_0$  be an f-table scheme in  $F$ . An *instance over  $f$*  is a function from coordinates over  $f$  to  $\text{DOM}(l_0)$ , which is defined over a finite set of coordinates. An *instance over  $S$*  is a function mapping each f-table  $f$  in  $F$  to an instance over  $f$ .

An *entry* of an f-table instance  $f$  is a coordinate over which  $f$  is defined. The actual value that  $f$  associates with an entry is called a *measure*. Note that measures and attributes are both defined with respect to levels of dimensions, and thus the distinction between them is terminological and not conceptual. In other words, our model does allow a symmetric treatment of measures and components of coordinates.

It is apparent that our notion of “symbolic coordinate” is related with that of “tuple” in the relational model. This is motivated by the intuition that an f-table is a “logical” counterpart of the “physical” notion of a multi-dimensional array. It can also be noted that the notation we use for symbolic coordinates resembles subscripting into a multi-dimensional array (although in a non-positional way). There is however an important difference between f-tables and multi-dimensional arrays. Specifically, in arrays, “physical” coordinates vary over intervals within linearly-ordered domains (in particular, over initial segments of natural numbers), whereas we do not pose any restrictive hypothesis on the domains over which coordinates range. In this sense, our notion of coordinate is “symbolic.”

**Example 8.** A possible instance for the sample scheme Toys defined in Example 5 is shown in Figure 2. Note that two different (graphical) representations for

f-tables are used in the figure. A symbolic coordinate over the f-table SALES is [day : Jan 5, 97, item : Scrabble, store : Navona]. The actual instance associates the measure 32 with this entry.

day	item	store	SALES
Jan 5, 97	Scrabble	Navona	32
Jan 5, 97	Risiko	Navona	27
Jan 5, 97	Lego	Sun City	42
Jan 5, 97	Risiko	Sun City	22
Feb 19, 97	Scrabble	Navona	32
Feb 19, 97	Lego	Navona	25
Feb 19, 97	Lego	Colosseum	11
Mar 10, 97	Risiko	Navona	5
Mar 10, 97	Lego	Sun City	6

PRICE-LIST	Jan-97	Feb-97	Mar-97
Lego	12. <sup>99</sup>	9. <sup>99</sup>	9. <sup>99</sup>
Risiko	14. <sup>99</sup>	12. <sup>99</sup>	12. <sup>99</sup>
Scrabble	12. <sup>99</sup>	12. <sup>99</sup>	12. <sup>49</sup>
Trivia		18. <sup>99</sup>	17. <sup>99</sup>

Fig. 2. A sample instance over the Toys scheme.

Figure 2 suggests that several different representations of a same f-table are possible. A tabular representation for an f-table  $f$  (like the one used for SALES) consists of a relation over the attributes of  $f$ , plus a further column for the measures provided by the instance; this representation suggests a way to implement f-tables with the relational model. If an f-table has  $n$  attributes, it can be also represented as a  $n$ -dimensional array (like the one used for PRICE-LIST) in which an entry corresponds to a measure of the instance. This representation recalls the way in which multidimensional systems usually store data.

### 3 The MultiDimensional Calculus by Examples

In this section, we present MD-CAL, a query language for the MD model. This language is a calculus for f-tables, and allows the analyst to express analytical queries in a declarative way.

Interpreted scalar and aggregate functions can be used in queries, but the semantics of the language is parametric with respect to them. This gives us the freedom of choosing the most suitable collection of functions, according to the specific application domain. Then, given a collection  $\mathcal{G}$  of interpreted functions, we denote by MD-CAL( $\mathcal{G}$ ) the MD query calculus that allows to use the functions in the collection  $\mathcal{G}$ .

The presentation is mainly based on examples that refer to the Toys sample scheme introduced in the previous section.

#### 3.1 Basic Queries

Intuitively, a *MultiDimensional query* is a mapping from instances over an input MD scheme to instances over an output MD scheme. The input and output

schemes are defined over the same dimensions but distinct f-tables. For the sake of simplicity, we shall assume that the output scheme of a query contains just a single f-table, called the *output f-table* of the query.

If the output f-table of a query has scheme  $f[A_1 : l_1, \dots, A_n : l_n] : l$ , then an MD-CAL *query* is specified by means of an expression of the following form.

$$\{x_1, \dots, x_n : x \mid \psi(x, x_1, \dots, x_n)\}$$

In the first part of the query, called the *target list*,  $x, x_1, \dots, x_n$  are distinct *variables*; the distinguished variable  $x$  is called the *result variable*. Furthermore,  $\psi(x, x_1, \dots, x_n)$  is a first-order formula in which  $x, x_1, \dots, x_n$  are the only free variables. The formula  $\psi$  is composed by equality atoms involving f-tables, roll-up functions, and interpreted scalar and aggregate functions.

Intuitively, the result of the query is an instance over the output f-table, associating a measure  $m$  to the entry  $[A_1 : c_1, \dots, A_n : c_n]$  for those values  $m, c_1, \dots, c_n$  that, respectively substituted to  $x, x_1, \dots, x_n$ , satisfy the formula.

An important aspect in MD-CAL is what we call “definiteness” of queries. Intuitively, this property guarantees that queries define indeed f-tables, which, by definition, must be finite and satisfy a sort of functional dependency from coordinates to measures. We shall discuss on this issue in Section 4.3; for the time being, we present only examples that obviously satisfy this property.

As a first example, the following query is used to define an f-table

$$\text{ROME-SALES}[day : day, item : item, store : store] : \text{numeric}$$

to represent the same information as SALES, but limited to the stores in Rome.

$$\{x_1, x_2, x_3 : x \mid x = \text{SALES}[day : x_1, item : x_2, store : x_3] \wedge \text{Rome} = \text{R-UP}_{\text{store}}^{\text{city}}(x_3)\}$$

### 3.2 Scalar Functions

As we have said, an atom in the formula of a query can use a predefined set  $\mathcal{G}$  of interpreted functions. This set can include system-defined or user-defined *scalar* functions, that is, functions that use only atomic values as inputs and outputs (e.g., all the standard mathematical operators, such as  $+$  and  $*$ ). Special care must be devoted in defining the semantics of a scalar function when one or more of its arguments is undefined. In what follows, unless otherwise stated, we will assume that the result of a function is undefined whenever one of its argument is undefined.

The following query defines the f-table with scheme

$$\text{DAILY-REVENUES}[day : day, item : item, store : store] : \text{numeric}$$

that represents the daily revenues, for each store and item. A measure for a certain item is obtained by multiplying the number of pieces sold in a day, by the price of the item in that month.



$$\{x_1, x_2, x_3 : x \mid \exists x_4 (x_4 = \text{R-UP}_{\text{day}}^{\text{month}}(x_1) \wedge x = \text{SALES}[\text{day} : x_1, \text{item} : x_2, \text{store} : x_3] * \text{PRICE-LIST}[\text{item} : x_2, \text{month} : x_4])\}$$

### 3.3 Aggregate Functions

The set  $\mathcal{G}$  can also include *aggregate* functions, that is, functions that applied to a collection of values yield an atomic value; these are of special interest in OLAP systems. Typical aggregate functions are those of SQL, that is, **min**, **max**, **count**, **sum**, and **avg**, which apply to expressions over columns.

For instance, the following query defines the f-table having scheme

$$\text{SUMMARY-SALES}[\text{week} : \text{week}, \text{item} : \text{item}, \text{area} : \text{area}] : \text{numeric},$$

which represents summary data of sales, detailed by week, item, and area.

$$\{x_1, x_2, x_3 : x \mid x = \text{sum}(y_1, y_2 : y \mid y = \text{SALES}[\text{day} : y_1, \text{item} : x_2, \text{store} : y_2] \wedge x_1 = \text{R-UP}_{\text{day}}^{\text{week}}(y_1) \wedge x_3 = \text{R-UP}_{\text{store}}^{\text{area}}(y_2))\}$$

The argument of the operator **sum** in the above query is a query  $q$  itself. The target list of  $q$  specifies “local” variables, and the result variable is used for the aggregation. Intuitively, the result of the whole query is as follows. For a week  $w$ , an item  $i$ , and an area  $a$ , let  $S_{w,i,a}$  be the result of the query obtained from  $q$  by substituting  $w, i, a$  for  $x_1, x_2, x_3$ , respectively. It is easy to see that  $S_{w,i,a}$  represents the number of sales of the item  $i$  in the days of the week  $w$ , in the stores of  $a$ . Now, let  $m$  the total sum of the sales in  $S_{w,i,a}$ . Then, the result of the whole query associates  $m$  with the coordinate  $w, i, a$ .

Note that, similarly to SQL, only entries of f-tables (which are non-null by definition) take part of the computation of the sum.

Assume now that we want to compute the f-table having scheme

$$\text{WEEKLY-REVENUES}[\text{week} : \text{week}, \text{item} : \text{item}, \text{store} : \text{store}] : \text{numeric},$$

which represents the weekly revenues, detailed by item and store. To do so, we can make use of the previously defined f-table DAILY-REVENUES, summarizing by weeks, as follows.

$$\{x_1, x_2, x_3 : x \mid x = \text{sum}(y_1 : y \mid y = \text{DAILY-REVENUES}[\text{day} : y_1, \text{item} : x_2, \text{store} : x_3] \wedge x_1 = \text{R-UP}_{\text{day}}^{\text{week}}(y_1))\}$$

However, we can also write the following query, which does not require the definition of DAILY-REVENUES.

$$\{x_1, x_2, x_3 : x \mid x = \text{sum}(y_1 : y \mid \exists y_2 (x_1 = \text{R-UP}_{\text{day}}^{\text{week}}(y_1) \wedge y_2 = \text{R-UP}_{\text{day}}^{\text{month}}(y_1) \wedge y = \text{SALES}[\text{day} : y_1, \text{item} : x_2, \text{store} : x_3] * \text{PRICE-LIST}[\text{item} : x_2, \text{month} : y_2]))\}$$

### 3.4 Abstraction Queries

In the context of multi-dimensional data, it is often useful to transform measures into components of coordinates of f-tables, and vice versa. We call *abstractions* such transformations. The following example shows how to perform an abstraction in MD-CAL. The query generates the boolean f-table

**TOTAL-SALES**[*item* : item, *store* : store, *year* : year, *items-sold* : numeric] : boolean

by summarizing on the number of sales and using the result as an element of a symbolic coordinate.

$$\{x_1, x_2, x_3, x_4 : x \mid x = \text{true} \wedge \\ x_4 = \text{sum}(y_1 : y \mid y = \text{SALES}[\text{day} : y_1, \text{item} : x_1, \text{store} : x_2] \wedge x_3 = \text{R-UP}_{\text{day}}^{\text{year}}(y_1))\}$$

Intuitively, the effect of this query is the following. For each triple  $c_1, c_2, c_3$  of values over the variables  $x_1, x_2, x_3$ , the two formulas in the body are evaluated, using  $x_4$  and  $x$  to hold the respective results, say,  $c_4$  and  $m$ . Then  $m$  (that is, true) is assigned to the entry having coordinate  $[c_1, c_2, c_3, c_4]$ . Note that the f-table we obtain represents, in every respect, a relation of the relational model.

## 4 The MultiDimensional Calculus

In this section we formally introduce the MultiDimensional query calculus MD-CAL.

In what follows we fix a MultiDimensional scheme  $\mathcal{S}$  and an instance  $\mathcal{I}$  over  $\mathcal{S}$ . We also fix a collection  $\mathcal{G}$  of scalar and aggregate interpreted functions. Each function in  $\mathcal{G}$  is characterized by a *signature* and an *interpretation*. For a *scalar function*  $g \in \mathcal{G}$ , the signature has the form  $g : l_1 \times \dots \times l_n \rightarrow l$ , where  $l, l_1, \dots, l_n$  are levels; an interpretation for  $g$  is a function from  $\text{DOM}(l_1) \times \dots \times \text{DOM}(l_n)$  to  $\text{DOM}(l)$ . For an *aggregate function*  $h \in \mathcal{G}$ , the signature has the form  $h : 2^{l'} \rightarrow l$ , where  $l$  and  $l'$  are levels; an interpretation for  $h$  is a function from finite *multi-sets* of elements in  $\text{DOM}(l')$  to elements of  $\text{DOM}(l)$ .

### 4.1 Syntax

For each level  $l$ , assume the existence of a countable set of *variables of type  $l$* .

The *terms* (over  $\mathcal{S}$  and  $\mathcal{G}$ ) and their respective types are recursively defined as follows.

- A variable of type  $l$  is a term of type  $l$ ;
- a value in  $\text{DOM}(l)$  is a term of type  $l$ ;
- if  $t$  is a term of type  $l$  and  $l$  rolls up to a level  $l'$ , then  $\text{R-UP}_l^{l'}(t)$  is a term of type  $l'$ ;
- if  $f[A_1 : l_1, \dots, A_n : l_n] : l$  is an f-table scheme and  $t_1, \dots, t_n$  are terms of type  $l_1, \dots, l_n$ , respectively, then  $f[A_1 : t_1, \dots, A_n : t_n]$  is a term of type  $l$ ;

- if  $g : l_1 \times \dots \times l_n \rightarrow l$  is a scalar function and  $t_1, \dots, t_n$  are terms of type  $l_1, \dots, l_n$ , respectively, then  $g(t_1, \dots, t_n)$  is a term of type  $l$ ;
- if  $h : 2^{l'} \rightarrow l$  is an aggregate function, and  $\tau \mid \psi$  is a query (defined below) whose result variable is of type  $l'$ , then  $h(\tau \mid \psi)$  is a term of type  $l$ .

An *atom* (over  $\mathcal{S}$  and  $\mathcal{G}$ ) is an expression of the form  $t = t'$ , where  $t$  and  $t'$  are terms (over  $\mathcal{S}$  and  $\mathcal{G}$ ) of the same type. The *formulas* (over  $\mathcal{S}$  and  $\mathcal{G}$ ) are defined as follows.

- An atom is a formula;
- if  $\psi_1$  and  $\psi_2$  are formulas, then  $\psi_1 \wedge \psi_2$ ,  $\psi_1 \vee \psi_2$ , and  $\neg \psi_2$  are formulas;
- if  $\psi$  is a formula and  $x$  is a variable, then  $\exists x(\psi)$  and  $\forall x(\psi)$  are formulas.

The notions of *free* and *bound* occurrences of variables are as usual, with the following additional consideration: the variables in the target list of an aggregation term are bound outside the term.

An MD-CAL *query* is an expression of the form

$$\{x_1, \dots, x_n : x \mid \psi(x, x_1, \dots, x_n)\},$$

where  $\psi(x, x_1, \dots, x_n)$  is a formula having  $x, x_1, \dots, x_n$  as distinct free variables. The expression  $x_1, \dots, x_n : x$  is called the *target list*, and  $x$  the *result variable*.

## 4.2 Semantics

Let  $q$  be an MD query of the form  $\{x_1, \dots, x_n : x \mid \psi(x, x_1, \dots, x_n)\}$ . The *pre-result* of  $q$  on  $\mathcal{I}$ , denoted by  $\text{PRE}(q(\mathcal{I}))$ , is the set of tuples of values  $c, c_1, \dots, c_n$  that, respectively substituted to  $x, x_1, \dots, x_n$ , satisfy the formula  $\psi$  with respect to  $\mathcal{I}$ . In such tuples, the first component  $c$  is called the *result value*.

The notion of *satisfaction* of a formula with respect to a substitution  $\theta$  and an instance  $\mathcal{I}$  is defined in the usual way, with the following considerations.

- The substitutions are typed, so that variables vary over values of the corresponding types. For the time being, we assume that values are chosen from the domain  $\text{DOM}(\mathcal{S})$ , that is, the union of the domains of the levels occurring in  $\mathcal{S}$ .
- Consider an atom of the form  $t = h(\tau \mid \psi)$ , where  $h$  is an aggregate function, and a substitution  $\theta$  over the free variables of the atom. Let  $T$  be the pre-result of the query  $\{\tau \mid \theta(\psi)\}$  over  $\mathcal{I}$  and let  $M$  be the multi-set containing the result values of  $T$ , with the respective multiplicity. Then, the atom is satisfied if  $\theta(t) = h(M)$ .

Thus, the pre-result of an MD-CAL query is a set of tuples, to be used as coordinates and measures of the result f-table. This is however not always possible, since there are pre-results that do not correspond to f-table instances. We say that the pre-result of a query over an instance is *functional* if it does not contain a pair of different tuples that coincide on all values, but the result value.

Let  $q$  be a query having  $f[A_1 : l_1, \dots, A_n : l_n] : l$  as output scheme. If the pre-result  $F = \text{PRE}(q(\mathcal{I}))$  of  $q$  is functional, then we can build in the natural way an f-table instance  $\text{FT}(F)$  from it, as follows. For each tuple  $c, c_1, \dots, c_n$  in  $F$ ,  $\text{FT}(F)$  associates the result value  $c$  to the symbolic coordinate  $[A_1 : c_1, \dots, A_n : c_n]$ . Then, the *result* of  $q$  over  $\mathcal{I}$ , denoted by  $q(\mathcal{I})$  is defined as  $\text{FT}(\text{PRE}(q(\mathcal{I})))$ .

### 4.3 Definiteness

Apart from functionality, the result of a query should satisfy another important property: the finiteness of the result. Actually, in the context of the relational calculus, a more general notion, the *domain independence*, has been defined to capture the finiteness of queries. In this section, we introduce and discuss the issue of *definiteness* as a desirable property for MD-CAL queries: intuitively, this notion combines the properties of domain independence (in the context of the MD model) and functionality.

Indeed, the notion of domain independence has been further generalized for queries involving interpreted functions, in particular, to *bounded depth domain independence* [1]. Now, it is straightforward to define the result of an MD-CAL query relativized to a domain  $D$  rather than to the domain  $\text{DOM}(\mathcal{S})$ . Then, we can say that, intuitively, an MD-CAL query  $q$  (using functions from a collection  $\mathcal{G}$ ) is bounded depth domain independent if, for any instance  $\mathcal{I}$ , its result depends only on a domain including the active domain  $\text{ADOM}(q, \mathcal{I})$  of  $\mathcal{I}$  and of  $q$ , plus a further small set of values obtained by applying a bounded number of times the roll-up functions and the functions in  $\mathcal{G}$  to  $\text{ADOM}(q, \mathcal{I})$ .

We say that an MD-CAL query  $q$  is *definite* if, for any input instance  $\mathcal{I}$ , it is bounded depth domain independent and functional.

Syntactic characterizations that ensure bounded depth domain independence have been proposed, for instance, *embedded allowedness* [10]. On the other hand, the property of functionality can be reduced to a problem of implication of *functional dependencies* for the MD-CAL language.

*Example 9.* Let us consider the query in Section 3.2 defining the f-table DAILY-REVENUES. Intuitively, this query is bounded depth domain independent since: (i) the variables  $x_1, x_2, x_3, x_4$  are bounded to values occurring in the input instance (note that the variable  $x_4$  is bounded also because its values can be obtained by applying a roll-up function to a bounded variable); and (ii) the variable  $x$  is bounded to values that can be obtained by a single application of the scalar function  $*$ . Moreover, the query is functional since the functional dependency  $x_1, x_2, x_3 \rightarrow x$  is implicated by the following facts: (i)  $x_4$  functionally depends on  $x_1$ , because of the roll-up function; and (ii)  $x$  functionally depends on  $x_1, x_2, x_3, x_4$ , because of the application of a scalar function to two measures that functionally depends on  $x_1, x_2, x_3$  and  $x_2, x_4$ , respectively. Hence, the query is definite.

If we restrict MD-CAL to queries involving no functions (neither roll-up nor interpreted ones), definiteness of MD-CAL queries corresponds to domain independence and functionality in the context of the relational model. It is well-know

that both properties are undecidable, but become decidable if the language is restricted to *positive existential* calculus queries [2]. It is also clear that definiteness is undecidable in MD-CAL, but decidable in positive existential MD-CAL without functions. We can show that definiteness is decidable for positive existential MD-CAL queries involving roll-up functions.

## 5 Expressive Power

In this section we study expressiveness of the MultiDimensional model and the MD-CAL query language. We show that the MD model subsumes the relational data model. We also show that, with suitable choices of interpreted functions, the conjunctive MD-CAL expresses the relational calculus (Section 5.1) and Klug's query languages with aggregate functions [16] (Section 5.2). In doing so, we show that a restricted number of functions suffices to express SQL with aggregation operators, and some of its extensions in the context of data analysis.

In what follows, given a data model  $m$ , we denote by  $\text{REP}_m(S)$  ( $\text{REP}_m(I)$ , respectively), the representation, in the MultiDimensional model, of the scheme  $S$  (instance  $I$ ) of the model  $m$ . Then, we say that an MD-CAL query  $q$  *expresses* a query  $q'$  in a language  $L$  for a model  $m$  if, for any instance  $I$  of  $m$  it is the case that  $\text{REP}_m(q'(I))$  equals  $q(\text{REP}_m(I))$ . We also say that an MD query language *expresses* another query language  $L$  if it expresses all the queries that are expressible in  $L$ .

### 5.1 MD and Relational Databases

Let  $S$  be a relational database scheme, that is, a set of relational schemes of the form  $R(A_1 : d_1, \dots, A_k : d_k)$ , where each  $A_i$  (with  $1 \leq i \leq k$ ) is an attribute name and each  $d_i$  is an associated domain. The *representation*  $\text{REP}_{rel}(S)$  of the scheme  $S$  is the MD scheme containing: (i) a dimension (consisting of a single level) for each distinct domain of  $S$ ; and (ii) an f-table scheme  $R[A_1 : l_1, \dots, A_k : l_k] : T$  for each relation scheme  $R(A_1 : d_1, \dots, A_k : d_k)$ , where each  $l_i$  is the level associated to the domain  $d_i$ , and  $T$  is a level whose domain contains only the boolean value "true."

Then, the *representation*  $\text{REP}_{rel}(I)$  contains, for each relation  $R \in S$ , an f-table instance  $R$  such that  $R[t] = \text{true}$  if and only if the tuple  $t$  belongs to  $R$  in the instance  $I$ .

Clearly, MD-CAL expresses the relational calculus and, therefore, the relational algebra. However, an interesting result is that the *conjunctive* MD-CAL language (involving only  $\exists$  and  $\wedge$ ) plus two simple interpreted functions is as expressive as the relational algebra. These functions are: (i) the aggregate function  $\Phi$ , which tests whether its argument is not empty; (ii) the scalar function **if**, to compose terms of the form **if**  $C$  **then**  $E$  **else**  $E'$ , whose first argument  $C$  is a conjunction of boolean terms, that returns the second argument  $E$  if  $C$  evaluates to "true," and the third argument  $E'$  if  $C$  evaluates to "false" or it is undefined.

**Theorem 10.** *Conjunctive MD-CAL( $\Phi$ , if) expresses the relational algebra.*

*Proof: (Sketch)* All the operators of the relational algebra, but the projection, the union, and the difference, can be easily implemented in conjunctive MD-CAL without interpreted functions. We then use the following expression to compute the result of a projection  $R = \pi_{A_1, \dots, A_k}(S)$ .

$$x = \Phi \left( \{x_1, \dots, x_k : x \mid y_{k+1}, \dots, y_n : y \mid y = S[A_1 : x_1, \dots, A_k : x_k, A_{k+1} : y_{k+1}, \dots, A_n : y_n]\} \right)$$

The difference of two relations  $R$  and  $S$  is expressed by:

$$x = \text{if } R[A_1 : x_1, \dots, A_n : x_n] \wedge S[A_1 : x_1, \dots, A_n : x_n] \\ \text{then } \perp \text{ else } R[A_1 : x_1, \dots, A_n : x_n],$$

where the symbol  $\perp$  stands for ‘undefined.’

Finally, the union of two relations  $R$  and  $S$  is expressed by:

$$x = \text{if } R[A_1 : x_1, \dots, A_n : x_n] \\ \text{then } R[A_1 : x_1, \dots, A_n : x_n] \text{ else } S[A_1 : x_1, \dots, A_n : x_n]$$

□

## 5.2 MD and Aggregate Functions

We now compare the expressive power of MD-CAL with the languages (an algebra and a calculus) having aggregate functions proposed by Klug [16]. In particular, Klug’s algebra, denoted by  $RA^{\text{Agg}}$ , is a standard relational algebra extended with an *aggregate formation* operator and a family **Agg** of aggregate functions. Intuitively, the aggregate formation operator partitions its argument according to a group-by list, and then applies an aggregate function  $h$  to each partition, to yield a tuple in the result. We assume that any aggregate function in **Agg** has a natural counterpart in the MD setting.

**Theorem 11.** *Conjunctive MD-CAL( $\Phi$ , if, **Agg**) expresses  $RA^{\text{Agg}}$ .*

Interestingly, there is a trade-off between scalar functions and the more complex aggregate functions. In particular, it turns out that any traditional aggregate function is subsumed by **sum** together with suitable scalar functions.

**Theorem 12.** *Conjunctive MD-CAL(**sum**, if,  $\geq$ ) expresses MD-CAL().*

Let  $\mathcal{G}_{sqi}$  be the set of SQL aggregation operators, that is, **sum**, **count**, **avg**, **min**, and **max**. We also have the following result.

**Theorem 13.** *Conjunctive MD-CAL(**sum**, if,  $+$ ,  $*$ ,  $/$ ,  $\geq$ ) expresses MD-CAL( $\mathcal{G}_{sqi}$ ).*

The intuition behind this result is that many statistical computations on numeric series are essentially based on the evaluation of the “expected value” of some scalar functions, defined as the mean of the function applied to the elements of the series. A similar result has been obtained in the context of a nested relational language [18].

Actually, we can show that conjunctive MD-CAL(**sum**, **if**, **+**, **\***, **/**, **≥**) expresses other aggregate functions, including the *generalized quantifiers* (set predicates such as **some**, **all**, and **at-least**) introduced in [19], and special operators (like **rank**, **ratio-to-report**, and **n-tile**) introduced in some SQL’s extensions in the context of data analysis [23].

*Example 14.* Let  $S[A : d] : \text{numeric}$  be an f-table that associates a numeric measure to elements of a dimension  $d$ . The following query defines the f-table  $S\text{-RANK}[A : d] : \text{numeric}$  that associates a rank with  $S$ . Specifically, if there are  $n$  distinct values that occur in the measure of  $S$ , then  $S\text{-RANK}$  associates a natural number with each entry of  $S$ :  $n$  with the entry having the highest value, and 1 with the entry having the lowest.

$$\{x_1 : x \mid x = \text{sum}(y_1 : y \mid y = \text{if } S[A : x_1] \geq S[A : y_1] \text{ then } 1 \text{ else } 0)\}$$

## 6 Conclusion

In this paper we have proposed a model and a calculus-based query language to establish a theoretical basis for multidimensional data.

Practical OLAP systems require a number of query languages, at different abstraction levels. On one hand, the final user should be enabled to perform point-and-click operations by means of graphical metaphors. Typical ways of manipulating a multi-dimensional data collection are: roll up (summarize data), drill down (go to more detailed data), slice and dice (select and project on a bidimensional view), pivot (reorient a data cube, projecting on different dimensions). On the other hand, the sophisticated user that needs to express more complex queries should be allowed to use a declarative, high-level language. Note that a practical language of this kind can be easily drawn from MD-CAL by adding some syntactical sugar. Finally, query optimization can be effectively performed by referring to a procedural, algebraic language. Thus, a family of different languages should be adopted by an OLAP system, and mapping between them should be defined.

Further current research topics in the context of OLAP systems are modeling and optimization. Dimensional modeling focuses on how information can be organized according to natural business concepts, i.e., the way decision-makers look at their business data, to enable decision support. Optimization concerns the ways in which factual data can be efficiently stored and manipulated. There are two main approaches to this problem in the context of decision-support applications: materialization of pre-computed summary data [13,21] and query optimization [3,7].

The formal nature of the model proposed here is well-suited for an investigation of the above problems. In particular, we are currently developing an algebra for the MultiDimensional model, for studying the efficient evaluation of multidimensional queries.

## Acknowledgments

We would like to thank Sophie Cluet and the anonymous referees for helpful suggestions.

## References

1. S. Abiteboul and C. Beeri. On the power of languages for the manipulation of complex objects. Technical Report 846, INRIA, 1988.
2. S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
3. S. Agarwal et al. On the computation of multidimensional aggregates. In *Twenty-second Int. Conf. on Very Large Data Bases, Bombay*, pages 506–521, 1996.
4. R. Agrawal, A. Gupta, and S. Sarawagi. Modeling multidimensional databases. In *Thirteenth IEEE International Conference on Data Engineering*, pages 232–243, 1997.
5. D. Chatziantoniou and K. Ross. Querying multiple features of groups in relational databases. In *Twenty-second Int. Conf. on Very Large Data Bases, Bombay*, pages 295–306, 1996.
6. S. Chaudhuri and U. Dayal. Decision support, Data Warehousing, and OLAP. In *Tutorials of the Twenty-second Int. Conf. on Very Large Data Bases*, 1996.
7. S. Chaudhuri and K. Shim. Optimization of queries with user-defined predicates. In *Twenty-second Int. Conf. on Very Large Data Bases, Bombay*, pages 87–98, 1996.
8. E.F. Codd, S.B. Codd, and C.T. Salley. Providing OLAP (On Line Analytical Processing) to user-analysts: An IT mandate. Arbor Software White Paper, <http://www.arborsoft.com>.
9. G. Colliat. OLAP, relational, and multidimensional database systems. *ACM SIGMOD Record*, 25(3):64–69, September 1996.
10. M. Escobar-Molano, R. Hull, and D. Jacobs. Safety and translation of calculus queries with scalar functions. In *Twelfth ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems*, pages 253–264, 1993.
11. J. Gray, A. Bosworth, A. Layman, and H. Pirahesh. Data Cube: a relational aggregation operator generalizing group-by, cross-tab, and sub-totals. In *Twelfth IEEE International Conference on Data Engineering*, pages 152–159, 1996.
12. M. Gyssens, L.V.S. Lakshmanan, and I.N. Subramanian. Tables as a paradigm for querying and restructuring. In *Fifteenth ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems*, pages 93–103, 1996.
13. V. Harinarayan, A. Rajaraman, and J. Ullman. Implementing data cubes efficiently. In *ACM SIGMOD International Conf. on Management of Data*, pages 205–216, 1996.
14. W.H. Inmon. *Building the Data Warehouse*. John Wiley, second edition, 1996.
15. T. Isakowitz, S. Schocken, and H.C. Lucas. Toward a logical/physical theory of spreadsheet modeling. *ACM Trans. on Inf. Syst.*, 13(1):1–37, January 1995.



16. A. Klug. Equivalence of relational algebra and relational calculus query languages having aggregate functions. *Journal of the ACM*, 29(3):699–717, 1982.
17. L. Libkin, R. Machlin, and L. Wong. A query language for multidimensional arrays: Design, implementation, and optimization techniques. In *ACM SIGMOD International Conf. on Management of Data*, pages 228–239, 1996.
18. L. Libkin and L. Wong. Aggregate functions, conservative extension, and linear orders. In *Workshop on Database Programming Languages*, pages 282–294, 1993.
19. S. Rao, A. Badia, and D. Van Gucht. Providing better support for a class of decision support queries. In *ACM SIGMOD International Conf. on Management of Data*, pages 217–227, 1996.
20. A. Shoshani. OLAP and statistical databases: Similarities and differences. In *Sixteenth ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems*, pages 185–196, 1997.
21. D. Srivastava, S. Dar, H.V. Jagadish, and A. Levy. Answering queries with aggregation using views. In *Twenty-second Int. Conf. on Very Large Data Bases, Bombay*, pages 318–329, 1996.
22. Stanford Technology Group, Inc. Designing the data warehouse on relational databases, 1995. Unpublished manuscript.
23. Red Brick Systems. Decision-makers, business data, and RSQL, 1995. White Paper, <http://www.redbrick.com>.
24. J.L. Weldon. Managing multidimensional data: Harnessing the power. *Database Programming & Design*, 8(8):24–33, August 1995.