

# A Database Array Algebra for Spatio-Temporal Data and Beyond

Peter Baumann

FORWISS (Bavarian Research Center for Knowledge-Based Systems)  
Orleansstr. 34, D-81667 Munich, Germany  
baumann@forwiss.de

**Abstract.** Recently multidimensional arrays have received considerable attention among the database community, applications ranging from GIS to OLAP. Work on the formalization of arrays frequently focuses on mapping sparse arrays to ROLAP schemata. Database modeling of further array types, such as image data, is done differently and with less rigid methods. A unifying formal framework for general array handling of image, sensor, statistics, and OLAP data is missing.

We present a cross-dimensional and application-independent algebra for the high-level treatment of arbitrary arrays. An array constructor, a generalized aggregate, plus a multidimensional sorter allow to declaratively manipulate arrays. This algebra forms the conceptual basis of a domain-independent array DBMS, RasDaMan, which offers an SQL-based query language with extensive algebraic query and storage optimization. The system is in practical use in neuro science.

We introduce the algebra and show how the operators transform to the array query language. The universality of our approach is demonstrated by a number of examples from imaging, statistics, and OLAP.

## 1 Introduction

In principle, any natural phenomenon becomes spatio-temporal array data of some specific dimensionality once it is sampled and quantised for storage and manipulation in a computer system; additionally, a variety of artificial sources such as simulators, image renderers, and data warehouse population tools generate array data. The common characteristic they all share is that a large set of large multidimensional arrays has to be maintained. We call such arrays *multidimensional discrete data* (MDD), expressing the variety of dimensions and separating them from the conceptually different multidimensional vectorial data appearing in geo databases.

As arrays obviously form both an important and a very clearly defined information category, it seems natural to describe them in a uniform manner through a homogeneous conceptual model. Preferably this is done in a way that the array model smoothly fits into existing overall models.

From a database perspective (and history), several separate information categories can be distinguished. Sets comprise the first category, well addressed by relational algebra and calculus. Semantic nets form the second one, being fundamentally different in structures and operations, although mappings to the relational model have been studied extensively. The third fundamental category is text, addressed by information retrieval (IR) technology. This distinction is not withstanding the fact that techniques to map object nets to relations with foreign keys are well-known, that IR techniques have found their way into relational products (e.g., Oracle8), and that hypertext combines nets and text into so-called semi-structured data. Arrays represent a separate fourth category, substantially different from the previous three. Again, mappings have been developed to the relational model, however, involving a significant semantic transformation. For sparse business data these are star, galaxy, and snowflake techniques [1], for image data these are blobs [2] (where a particularly high loss of semantics is incurred). A clear indicator for the semantic mismatch of SQL-based multidimensional queries is the resulting lack of functionality and performance, leading to several suggestions for extending the relational model – e.g., [3, 4, 5].

Multidimensional database research has history, as statistical databases have been studied since long [6, 7]; more recently, OLAP continues this tradition with a strong focus on business data [8]. Several proposals exist to formalize array structures and operations for OLAP [4, 3, 9, 10, 11, 5], for scientific computing [12] and for imaging [13]. The Discrete Fourier Transform (DFT) has been studied from a database viewpoint [14]. Often, however, formal concepts have not been implemented in an operational system and they have not been evaluated in real-life applications. Moreover, many of the formal models have been designed specifically with the goal of mapping arrays to relation tuples and in a way that, in practice, makes sense only for sparse arrays.

In this paper, we propose an algebraic framework (see [15] for a first version) which allows to express cross-dimensional queries, i.e., operations on arrays of any number of dimensions, simultaneously in one and the same expression and symmetric in all dimensions. Essentially, this algebra consists of only three operations: an array constructor, a generalized aggregation, and a multidimensional sorter. This core model does not rely on recursion and is safe in evaluation, yet it is sufficient to express a wide range of imaging, statistical, and OLAP operations. Therefore, our algebra can be seen as a “universal” framework, independent from the particular application domain. The concepts are implemented in the domain-independent array DBMS RasDaMan<sup>1</sup> [16, 15, 17], hence the name *RasDaMan Array Algebra*, or short: *Array Algebra*.

The remainder of this paper is organized as follows. In Section 2, Array Algebra is presented, together with practical examples from diverse application fields to illustrate its applicability. The step to an SQL-embedded array query language, RasQL, is shown in Section 3. Section 4 surveys related work, and Section 5 summarizes our findings.

---

<sup>1</sup> Raster Data Manager; see [www.forwiss.de/~rasdaman](http://www.forwiss.de/~rasdaman)

## 2 Array Algebra

We treat arrays as functions mapping n-dimensional points (i.e., vectors) from discrete Euclidean space to values. This is common in imaging for a long time – see, e.g., [18] – and has been transposed to database terminology in [16, 15]. To smoothly embed Array Algebra into existing overall algebras we use a set-oriented basis. Due to space constraints we have to omit most proofs here.

Operations on such arrays frequently apply a function simultaneously to a set of cells, requiring second-order functionals in the algebra. In practice they are necessary to allow for binding variables to points for iterating coordinate sets and also to aggregate arrays (or part thereof) into scalar values. The latter operation corresponds very much to relational set aggregators; however, instead of providing a limited list of aggregation operations as in the relational algebra, a general constructor is introduced by Array Algebra which is parametrized with the underlying base operation.

### 2.1 N-Dimensional Interval Arithmetics

We first introduce some notation for n-dimensional integer interval arithmetics. We call the coordinate set of an array its *spatial domain*. Informally, a spatial domain is defined as a set of n-dimensional points (i.e., algebraic vectors) in Euclidean space forming a finite hypercube with boundaries parallel to the coordinate system axes.

We assume common vector notation. For a natural number  $d > 0$ , we write  $x = (x_1, \dots, x_d) \in X \subseteq \mathbf{Z}^d$  for some d-dimensional vector  $x$ ,  $x+y$  for vector addition, etc. The point set forming the geometric extent of an array is called its *spatial domain*. A spatial domain  $X$  of dimension  $d$  spanned by  $\underline{l}$  and  $\underline{h}$  is defined as

$$\begin{aligned} X = [l_1 : h_1, \dots, l_d : h_d] &:= \bigcap_{i=1}^d \{x_i : l_i \leq x_i \leq h_i\} \text{ if } \forall 1 \leq i \leq d: l_i \leq h_i \\ &:= \{\} \text{ otherwise.} \end{aligned}$$

Functions  $lo, hi: \mathbf{P}(\mathbf{Z}^d) \rightarrow \mathbf{Z}^d$  (where  $\mathbf{P}$  is the Powerset) defined as  $lo(X) = \underline{l}$  and  $hi(X) = \underline{h}$  for some spatial domain  $X$  given as before denote the bounding vectors. We will abbreviate  $lo_i(X) = l_i$  and  $hi_i(X) = h_i$  for the  $i^{\text{th}}$  component. Function  $\dim(X) = d$  denotes the dimension of spatial domain  $X$ .

On such hypercubes, point set operations can be defined in a straightforward way. We admit only those operations which respect closure, such as `intersect` and `union*`, whereby the asterisk “\*” denotes the hull operation applied to the result:

$$\begin{aligned} \text{intersect}^*(X, Y) &:= \\ &[ \max(lo_1(X), lo_1(Y)) : \min(hi_1(X), hi_1(Y)), \dots, \\ &\quad \max(lo_d(X), lo_d(Y)) : \min(hi_d(X), hi_d(Y)) ] \end{aligned}$$

$$\begin{aligned} \text{union}^*(X, Y) &:= \\ &[ \min(\text{low}_1(X), \text{low}_1(Y)) : \max(\text{hi}_1(X), \text{hi}_1(Y)), \dots, \\ &\quad \min(\text{low}_d(X), \text{low}_d(Y)) : \max(\text{hi}_d(X), \text{hi}_d(Y)) ] \end{aligned}$$

Obviously these operations are commutative, associative, and distributive.

The `shift` operator allows to change a spatial domain's position according to a translation vector  $\underline{t}$ :

$$\text{shift}_{\underline{t}}(X) := \{ \underline{x} + \underline{t} : \underline{x} \in X \}$$

Let  $X$  be spanned by  $d$ -dimensional vectors  $\underline{l}$  and  $\underline{h}$ . For some integer  $i$  with  $1 \leq i \leq d$  and a one-D integer interval  $\mathbb{I} = [m:n]$  with  $l_i \leq m \leq n \leq h_i$ , the *trim of  $X$  to  $\mathbb{I}$  in dimension  $d$*  is defined as

$$\tau_{i, \mathbb{I}}(X) := \{ \underline{x} \in X : m \leq x_i \leq n \} = [ l_1 : h_1, \dots, m : n, \dots, l_d : h_d ]$$

Intuitively speaking, trimming slices off those parts of an array which are lower than  $m$  and higher than  $n$  in the dimension indicated; the dimension is unchanged. As opposed to this, a section cuts out a hyperplane with dimension reduced by 1. Formally, for some  $X$  as above, an integer  $p$  with  $1 \leq p \leq d$ , the *section of  $X$  at position  $p$  in dimension  $i$*  is given by

$$\begin{aligned} \sigma_{i, p}(X) &:= \{ \underline{x} \in \mathbb{Z}^{d-1} : \underline{x} = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_d), \\ &\quad (x_1, \dots, x_{i-1}, p, x_{i+1}, \dots, x_d) \in X \} \\ &= [ l_1 : h_1, \dots, l_{i-1} : h_{i-1}, l_{i+1} : h_{i+1}, \dots, l_d : h_d ] \end{aligned}$$

Trimming is commutative and associative, whereas a section changes dimension numbering and, therefore, has neither of these properties.

## 2.2 The Core Algebra

Let  $X \subseteq \mathbb{Z}^d$  be a spatial domain and  $F$  a homogeneous algebra. Then, an  $F$ -valued  $d$ -dimensional array over spatial domain  $X$  – or short: (*multidimensional*) array – is defined as

$$a : X \rightarrow F \text{ (i.e., } a \in F^X), \quad a = \{ (\underline{x}, a(\underline{x})) : \underline{x} \in X, a(\underline{x}) \in F \}$$

Array elements  $a(\underline{x})$  are referred to as *cells*. For notational convenience, we also allow to enumerate the components of a cell coordinate vector, e.g.,  $a(x_1, x_2, x_3)$ . Auxiliary function  $\text{sdom}(a)$  denotes the spatial domain of some array  $a$ ; further, we lift function  $\text{dim}$  to arrays. For an array  $a : X \rightarrow F$ ,  $\text{sdom}$  and  $\text{dim}$  are defined as

$$\begin{aligned} \text{sdom}(a) &:= X \\ \text{dim}(a) &:= \text{dim}(\text{sdom}(a)) \end{aligned}$$

The  $i^{\text{th}}$  dimension range of an array's spatial domain we will denote by  $\text{sdom}_i(a)$ .

**Example:** For a  $1024 \times 768$  image  $a$  with lower left corner in the origin of the coordinate system,  $\text{sdom}(a)=[0:1023,0:767]$ ,  $\text{dim}(a)=2$ .

The first functional we introduce is the *array constructor* MARRAY. It allows to define arrays by indicating a spatial domain and an expression which is evaluated for each cell position of the spatial domain. An iteration variable bound to a spatial domain is available in the cell expression so that a cell's value can depend on its position. Let  $X$  be a spatial domain,  $F$  a value set, and  $v$  a free identifier. Let further  $e_v$  be an expression with result type  $F$  containing zero or more free occurrences of  $v$  as placeholder(s) for an expression with result type  $X$ . Then, an *array over spatial domain  $X$  with base type  $F$*  is constructed through

$$\text{MARRAY}_{X,v}(e_v) = \{ (\underline{x}, a(\underline{x})) : a(\underline{x}) = e_{\underline{x}}, \underline{x} \in X \}$$

**Example:** Consider scaling of a greyscale image  $a$  with  $\text{sdom}(a) = [1:m, 1:n]$  by a factor  $s \in \mathbf{R}$ . We assume componentwise scalar division and rounding on vectors and write

$$\text{MARRAY}_{[1:m*s, 1:n*s], v}(a(\text{round}(v/s)))$$

For  $0 < s < 1$  the image is sized down; the interpolation method then corresponds to “nearest neighbor”, the simplest interpolation technique used in imaging.

The operation which in some sense is the dual to the MARRAY constructor is the *condenser* COND. It takes the values of an array's cells and combines them through the operation provided, thereby obtaining a scalar value. Again, an iterator variable is bound to a spatial domain to address cell values in the condensing expression. Let  $\circ$  be a commutative and associative operation with signature  $\circ: F, F \rightarrow F$ , let further  $v$  be a free identifier,  $X = \{ \underline{x}_1, \dots, \underline{x}_n \mid \underline{x}_i \in \mathbf{Z}^d \}$  a spatial domain consisting of  $n$  points, and  $e_{a,v}$  an expression of result type  $F$  containing occurrences of an array  $a$  and identifier  $v$ . Then, the *condense of  $a$  by  $\circ$*  is defined as

$$\text{COND}_{\circ, X, v}(e_{a,v}) := \bigcirc_{\underline{x} \in X} e_{a, \underline{x}} = e_{a, \underline{x}_1} \circ \dots \circ e_{a, \underline{x}_n}$$

**Examples:** Let  $a$  be the image as defined in the above example. Average pixel intensity is given by

$$\text{COND}_{+, \text{sdom}(a), v}(a(v)) / |\text{sdom}(a)| = \sum_{\underline{x} \in [1:m, 1:n]} a[\underline{x}] / (m*n)$$

For color table computation needed, e.g., for generation of a GIF image encoding, one has to know the set of all values occurring in array  $a$ . The condenser allows to derive this set by performing the union of all cell values:

$$\text{COND}_{\cup, \text{sdom}(a), v}(\{ a(v) \})$$

The third and last operator is an array sorter which proceeds along a selected dimension to reorder the corresponding hyperslices. Function  $\text{sort}_s$  rearranges a given array along a specified dimension  $s$  without changing its value set or spatial domain. To this end, an order-generating function is provided which associates a “sequence position” to each  $(d-1)$ -dimensional hyperslice. Let  $a$  be a  $d$ -dimensional array,  $i \in \mathbb{N}$  with  $1 \leq i \leq d$  a dimension number, and  $f_{s,a}: \text{sdom}_s(a) \rightarrow \mathbb{N}$  a total function which, for a given array  $a$ , inspects  $a$  in the sorting dimension  $s$  and delivers an ordering measure for each hyperslice. Further, let  $\text{perm}(\underline{x}, \underline{y})$  be a predicate indicating that vector  $\underline{x}$  is a permutation of vector  $\underline{y}$  (and vice versa). Then, the two sorters  $\text{sort}_{s,f}^{\text{asc}}$  and  $\text{sort}_{s,f}^{\text{desc}}$  for ascending and descending order, resp., are given as those arrays which consist of permutations of the hyperslices in the sort dimension and, additionally, fulfil the sorting criterion given by  $f$ :

$$\begin{aligned} \text{sort}_{s,f}^{\text{asc}}(a) := & \\ \{ (\underline{y}, b(\underline{y})) : \underline{y} \in \text{sdom}_s(a), & \\ \forall p, q \in \text{sdom}_s(a) : p < q \Rightarrow f_{s,b}(p) \leq f_{s,b}(q), & \\ \text{perm}((b(x_1, \dots, x_{s-1}, \text{sdom}_s(a).lo, x_{s+1}, \dots, x_d), \dots, & \\ b(x_1, \dots, x_{s-1}, \text{sdom}_s(a).hi, x_{s+1}, \dots, x_d)), & \\ (a(x_1, \dots, x_{s-1}, \text{sdom}_s(a).lo, x_{s+1}, \dots, x_d), \dots, & \\ a(x_1, \dots, x_{s-1}, \text{sdom}_s(a).hi, x_{s+1}, \dots, x_d))) & \} \end{aligned}$$

$$\begin{aligned} \text{sort}_{s,f}^{\text{desc}}(a) := & \\ \{ (\underline{y}, b(\underline{y})) : \underline{y} \in \text{sdom}_s(a), & \\ \forall p, q \in \text{sdom}_s(a) : p < q \Rightarrow f_{s,b}(p) \geq f_{s,b}(q), & \\ \text{perm}((b(x_1, \dots, x_{s-1}, \text{sdom}_s(a).lo, x_{s+1}, \dots, x_d), \dots, & \\ b(x_1, \dots, x_{s-1}, \text{sdom}_s(a).hi, x_{s+1}, \dots, x_d)), & \\ (a(x_1, \dots, x_{s-1}, \text{sdom}_s(a).lo, x_{s+1}, \dots, x_d), \dots, & \\ a(x_1, \dots, x_{s-1}, \text{sdom}_s(a).hi, x_{s+1}, \dots, x_d))) & \} \end{aligned}$$

The resulting array has the same number of dimensions, spatial domain, and base type as the input array. Note that function  $f_{s,a}$  has all degrees of freedom to assess any of  $a$ 's cell values for determining the measure value of a hyperslice on hand - it can be a particular cell value in the current hyperslice, the average of all hyperslice values, or even neighbored slices (e.g., for relative increases of sales values).

**Example:** Let  $a$  be a 1-D array with spatial domain  $D=[1:d]$  where cell values denote sales figures over time. Let further sorting function  $f_{s,a}$  be given as  $f_{s,a}(p) = a[p]$ . Then,  $\text{sort}_{0,f}^{\text{desc}}(a)$  delivers the ranked sales.

As an aside we note that the sort operator includes the relational *group by*. Below we will demonstrate that slice and roll-up operations arising from array access based on dimension hierarchies can be expressed, although - not very comfortably - by indicating the cell coordinates pertaining to a particular member set. Concepts for an intuitive, symbolic treatment of dimension hierarchies are currently under investigation.

## 2.3 Derived Operators

Several useful operations can be derived from the above ones. We present a selection of those which have turned out particularly important in practical applications.

### 2.3.1 Trimming and Section

The previously introduced spatial domain operations trimming and section give rise to corresponding array operations. For some array  $a$ , an 1-D interval  $I$ , and two natural numbers  $1 \leq t \leq \dim(a)$  and  $p \in \text{sdom}_d(a)$  they are defined as

$$\begin{aligned} \text{TRIM}_{t,I}(a) &:= \text{MARRAY}_{X,V}(a(v)) \text{ for } X=\tau_{t,I}(\text{sdom}(a)) \text{ and } d < \dim(a) \\ \text{SECT}_{t,p}(a) &:= \text{MARRAY}_{X,V}(a(v)) \text{ for } X=\sigma_{t,p}(\text{sdom}(a)) \text{ and } d < \dim(a) \end{aligned}$$

**Example:** Slicing of an OLAP cube  $c$  with spatial domain  $\text{sdom}(c)=D \times R \times P$  to extract subcube  $D' \times R' \times P' \subseteq \text{sdom}(c)$  is denoted as

$$\text{TRIM}_{1,D'}(\text{TRIM}_{2,R'}(\text{TRIM}_{3,P'}(c)))$$

### 2.3.2 Induced Operations

A basic set of operations is induced by the algebra of the underlying value sets. If  $a, b \in F^x$  are arrays and  $o$  is a binary operation on  $F$ , then  $o$  induces a binary operation on  $F^x$  denoted by  $o_{\text{ind}}$  such that, if  $c = a o_{\text{ind}} b$ , then  $c \in F^x$  and, for all  $x \in X$ ,  $c(x) = a(x) o b(x)$ . Along this line, we also allow to induce unary operations. Notably, these operations are not axiomatic; for a unary function  $f$  and a binary function  $g$ ,

$$\begin{aligned} \text{IND}_f(a) &= \text{MARRAY}_{X,V}(f(a(v))) \text{ for } X=\text{sdom}(a) \\ \text{IND}_g(a,b) &= \text{MARRAY}_{X,V}(g(a(v),b(v))) \text{ for } X=\text{sdom}(a)=\text{sdom}(b) \end{aligned}$$

Algebraic properties of  $F$  transform to corresponding structures on the set  $F^x$  of induced functions. If  $F$  is a field, then  $F^x$  is a vector space; for a ring  $F$ ,  $F^x$  is a module for suitably defined spatial domains.

**Examples:** Let  $a$  be a grayscale image over spatial domain  $X$ . Increasing intensity by 5 can be accomplished through induction on unary "+5":

$$\text{IND}_{+5}(a) = \{ (\underline{x}, b(\underline{x})) : b(\underline{x}) = a(\underline{x}) + 5, \underline{x} \in X \}$$

Consider now another grayscale image  $b$  over the same spatial domain  $X$ . Then, pixel addition can be induced to obtain image addition:

$$\text{IND}_{+}(a, b) = \{ (\underline{x}, c(\underline{x})) : c(\underline{x}) = a(\underline{x}) + b(\underline{x}), \underline{x} \in X \}$$

When used in a query, binary induction obviously implies a spatial join.

Let now  $c$  be a color image where the cell type is a three-integer record of red, green, and blue intensity, resp. Such a *pixel-interleaved* image is transformed into a *channel-interleaved* representation, i.e., three separate color planes, through induction on the record access operator ".", obtaining

$$\langle c.\text{red}, c.\text{green}, c.\text{blue} \rangle$$

The above type of induction is also referred to as pointwise induction, as points pairwise match for each application of the base function.

### 2.3.3 Aggregation

Obviously, the condenser provides the appropriate basis for aggregation over arrays. Table 1 lists some of the most common aggregations and their definition in Array Algebra.

**Table 1:** Some possible aggregate operators on arrays. Assumed are array expressions  $a$  (without restriction),  $b$  of result type Boolean, and  $c$  with a numerical result type.

Array aggregate definition	Meaning
Count_cells( $a$ )	The number of cells in $a$
Some_cells( $b$ )	Is there any cell in $b$ with value true?
All_cells( $b$ )	Do all cells of $b$ have value true?
Sum_cells( $c$ )	The sum of all cells in $c$
Avg_cells( $c$ )	The average of all cells in $c$
Max_cells( $c$ )	The maximum of all cells in $c$

## 2.4 Further Application Examples

A basic requirement in the development of Array Algebra has been to cover all applications of arrays in databases, the most important ones being statistics, OLAP, and imaging. To illustrate applicability of Array Algebra to these areas, we now present some advanced examples.



### 2.4.1 Statistics

**Example matrix multiplication:** Let  $a$  be an  $m \times n$  matrix and  $b$  an  $n \times p$  matrix. Then, the  $m \times p$  matrix product

$$a * b = \sum_{j=1}^n a_{i,j} * b_{j,k}$$

in Array Algebra is expressed as

$$\text{MARRAY}_{[1:m, 1:p], (i, k)} ( \text{COND}_{+, [1:n], j} ( a(i, j) * b(j, k) ) )$$

**Example auto correlation:** For two observation vectors  $x$  and  $y$  of dimension  $n$ , empirical covariance  $m_{x, y}$  is defined as

$$m_{x, y} = \frac{1}{n-1} \sum_{i=1}^n (x_i - x_{avg})(y_i - y_{avg})$$

In Array Algebra, the mean is given by  $x_{avg} = \text{COND}_{+, [1:n], i} (x(i)) / n$  and  $y_{avg} = \text{COND}_{+, [1:n], i} (y(i)) / n$ . Then,  $m_{x, y}$  is described in a straightforward manner:

$$\text{COND}_{+, [1:n], i} ( (x(i) - x_{avg}) * (y(i) - y_{avg}) ) / (n-1)$$

**Example histogram:** A histogram contains, for each possible value, the number of cells conveying this value. For some  $n$ -D one-byte integer array with intensity values between 0 and 255, the histogram is computed as

$$\text{MARRAY}_{[0:255], n} ( \text{COND}_{+, \text{sdom}(a), v} ( \text{if } a(v)=n \text{ then } 1 \text{ else } 0 \text{ fi} ) )$$

As we can see, the combination of MARRAY and COND appears in quite different contexts. Indeed, this type of operation forms the basis for an extremely wide range of analysis functions, such as statistical analyses, advanced OLAP consolidation operations like roll-up, slice&dice, as well as scaling, convolutions and filtering in image processing. It is capable of completely changing dimensionality, size, and cell types of arrays.

### 2.4.2 OLAP

**Example roll-up:** Let  $c$  be a sales datacube with  $D \times R \times P = [1:d, 1:r, 1:p]$  as spatial domain where dimension 1 counts days from 1 to today, dimension 2 enumerates sales regions, and dimension 3 contains products sold; cell values shall represent sales figures. The weekly average of sales per product and region, then, is expressed as

$$\text{MARRAY}_{[1:\text{today}/7]\times R\times P, (w,r,p)} ( \text{COND}_{+, [0:6], d} ( c(7*w+d, r, p) / 7 ) )$$

Aggregating over all products leads only to a slight change in the expression

$$\text{MARRAY}_{[1:d/7]\times R, (w,r)} ( \text{COND}_{+, [0:6]\times P, (d,p)} ( c(7*w+d, r, p) / 7 ) )$$

Notably, such queries can be of considerable length when formulated relationally, and usually involve several joins.

**Example top performers:** On the same cube, the top performing weeks are determined as follows. We use the notation  $\langle s:\text{sales}, w:\text{week} \rangle$  to describe a two-component record with component names *sales* and *week*. With function  $f_c$  given as  $f_c(i) = c[i].\text{sales}$ , the following expression rolls up this cube from days to weeks and delivers the accumulated sales over all regions and products of the top 3 weeks:

$$\text{sort}_{0, f_c^{\text{desc}}} ( \text{MARRAY}_{[1:d/7], w} ( \langle \text{COND}_{+, [0:6]\times R\times P, (d,r,p)} ( c[7*w+d, r, p] ) : \text{sales}, w:\text{week} \rangle ) ) [1:3].\text{week}$$

The last query heavily makes use of the fact that coordinate and cell values (dimension and measure elements in OLAP terminology) can be used interchangeably.

### 2.4.3 Imaging

**Example skewed section:** A skewed section through a 2-D image where the cutting line is not axis-parallel (Fig. 1) can be described by placing a skew factor  $s > 1$  on the indexing point  $\underline{x}$ , resulting in the shifted point position  $(s*\underline{x}_1, \underline{x}_1)$ :

$$\text{MARRAY}_{\text{sdom2}(a), v} ( a(s*v_1, v_1) )$$

Using the contents of another (1-D) array for indexing allows to pick arbitrary cells. Let  $a$  have  $\text{sdom}(a) = [1:m, 1:n]$  and  $s$  be a 1-D array with spatial domain  $X = [1:n]$ . Cell values of array  $s$  are used to index  $a$  (Fig. 1):

$$\text{MARRAY}_{\text{sdom1}(a), x} ( a(s(x), x) )$$

**Example filtering:** The following expression, parametrized over array  $a$  and mask  $m$  (such as the edge detector illustrated in Fig. 2) can be used as a template for general filtering operations:

1	3	1
0	0	0
-1	-3	-1

Through instantiation with mask  $s_1$  as given by Fig. 1 and mask  $s_2$  as the transpose of  $s_1$ , we can express the Sobel edge detector (see Fig. 2 for mask definition, Fig. 3 for an application example):

**Fig. 3:** Sobel filter applied to a 2-D raster image.

We observe that in many cases operations can be formulated without explicitly referring to the array dimension, allowing to develop parametrized cross-dimensional, domain-independent query libraries which go well beyond the capabilities of object-relational ADTs [19].

### 3 From Array Algebra to RasQL

In this Section, we sketch how Array Algebra translates to the query language RasQL. Array Algebra has been developed in the course of implementing this fully-fledged, domain-independent array DBMS based on an SQL-based query language, obeying strict data independence. Arrays are embedded as a data abstraction allowing, e.g., to define array-valued object or tuple attributes, depending on the hosting data model.

Arrays can be defined either concisely with dimension and extent per dimension fixed, or with a fixed number of dimensions but free lower or upper bounds in some dimension(s), or with dimension and boundaries left completely open. Runtime range checking on instances, then, is performed according to the amount of information provided in the data dictionary.

RasDaMan commits itself to the ODMG standard [20], hence the RasDaMan query language RasQL also follows the flavour of ODMG’s OQL which, in turn, leans itself on standard SQL-92. Queries range over collections which contain the class extents. Array expressions can appear both in the `select` and in the `where` clause of a query. The MARRAY equivalent in RasQL has the structure

```
marray <iterator> in <spatial domain>
values <expression>
```

The COND statement is somewhat extended. In the syntactic structure

```
condense <op>
over      <iterator> in <spatial domain>
where     <condition>
using     <expression>
```

The `where` condition allows to further restrict the cell set inspected. This makes thresholding and similar tasks more elegant to phrase and, in particular, supports optimization.

Optimizing MARRAY and COND expressions is not easy (although not impossible) due to the generality of the operators. We therefore continuously investigate on special cases where particularly efficient solutions exist; a rich set of over 100 rules has been identified and implemented yet [21]. For optimizability reasons and due to their practical importance, trimming, section, and induction are supported by special constructs; likewise, the condenser specializations mentioned in Table 1 are supported directly. Table 2 demonstrates some of these constructs with the help of application examples.

**Table 2:** Sample RasQL queries. We assume array-valued attributes for 2-D Landsat satellite images, 3-D volumetric images, and 3-D OLAP cubes, to be embedded in object classes (or relations, resp.).

Algebra operator	RasQL example	Explanation
$IND_f$	Select img + 5 from LandsatImages as img	The red channel of all Landsat images, intensified by 5
$IND_g$	Select oid(br) from BrainImages as br, BrainAreas as mask where br * mask > t	OIDs of all brain images where, in the masked area, intensity exceeds threshold value t
TRIM	Select w[ ***, y0:y1, z0:z1 ] from Warehouse as w	OLAP slicing (“***” exhausts the dimension)
SECT	Select v[ x, ***, *** ] from VolumetricImages as v	A vertical cut through all volumetric images
MARRAY	Select marray n in [0:255] values condense + over x in sdom(v) using v[x]=n from VolumetricImages as v	For each 3-D image its histogram <sup>2</sup>
COND	Select condense + over x in sdom(w) using w[x] > t from Warehouse as w	For each datacube in the warehouse, count all cells exceeding threshold value t

A trim expression in the left-hand side of an update assignment indicates the array part to be updated:

```
update <collection>
set   <array attribute>[<trim expression>]
assign <array expression>
where ...
```

Besides updating part of an array, this statement can also be used to extend an array by appending data (e.g., during periodical warehouse population or slicewise insertion

<sup>2</sup> RasQL supports the interpretation of Boolean values as numerics as is usual in many programming languages.

into a 3-D image), provided the affected dimension has been defined variable in the attribute definition. Formally the process of extending an array  $a$  in direction  $i$  with another array  $b$  matching  $a$  in all dimensions (with a possible exception in dimension  $i$ ) and base type is governed by the algebra expression as follows.

Let  $\underline{t}=(0,\dots,0,\text{sdom}_i(a),0,\dots,0)$  be the translation vector consisting of zeros except for component  $i \leq d$ . Then,

```

extend(a, b, i) :=
  MARRAYsdom(a) ∪ shift $\underline{t}$ (sdom(b)), v (
    if  $v \in \text{sdom}(a)$  then  $a(v)$  else  $b(v - \underline{t})$  fi
  )

```

## 4 Related Work

In this Section we survey related work in formalization of arrays in databases; in part we rely on the classification published in [22] which gives particular attention to the independence of array formalisms from their implementation.

Let us start, however, with APL [23,24] as a prominent representative from the programming languages area, dedicated to  $n$ -dimensional array manipulation. APL basically has to be compared not with the algebra but with the query language RasQL. On this level, both share interesting implementation problems which, however, are out of scope here. Conceptually, Array Algebra functionality has its respective counterparts in APL: the *enclose* operator " $\llcorner$ " corresponds to the MARRAY constructor; the *reduction* operator  $\diagup$  corresponds to the condenser, but is applied only to the outermost dimension; the *each* operator  $\lrcorner$  and, to some extent, *scalar functions* correspond to unary and binary induction. As a programming language with procedural constructs and recursion, APL is more powerful than Array Algebra, but not safe.

In the database world, the algebra underlying the EXTRA/EXCESS database system supports 1-D, variable-length arrays [25]. As for the operators, there is a function SUBARR corresponding to the Array Algebra operator SECT, and ARR\_APPLY corresponding to our unary induce  $\text{IND}_F$ . Aggregation is not supported.

Gray et al. [4] propose an SQL cube operator which generalizes *group-by*. It is based on a particular mapping of sparse arrays to relations. There is no clear separation between conceptual (multidimensional) model and the proposed SQL extensions; specifically, no formal algebra is provided.

In [3] a formal model for sparse array maintenance in relational systems (ROLAP) is presented. Array data are organized into one or more hypercubes whereby a cell value can either be an  $n$ -tuple (i.e., one nesting level of record elements) or a Boolean value denoting existence of the respective value combination. The algebra consists of a set of basic operations which are parametrized by user-defined functions. For example, operations pull and push increase/decrease, resp., a cube's dimension by changing coordinate values to cell contents and vice versa; in our example *top performers* we demonstrate how this is done in Array Algebra. Further, there is a join

operation to combine two arrays sharing  $k$  dimensions. The “join partners” are specified through user-defined functions outside the formalism. The same way aggregation is handled through functions outside the formalism as opposed to Array Algebra where the condenser serves to describe aggregation.

Cabibo and Torlone [9] propose a more “cube-oriented” formal multidimensional model and a corresponding query language based on a logical calculus. The data model relies on the notion of  $n$ -dimensional *f-tables*, i.e., (mathematical) relations where each cell is represented by a tuple of  $n$  coordinates and the cell value itself, which must be atomic. Aside of the usual logical quantifiers and connectors there are scalar and aggregation functions which are user-defined, hence outside the formalism. The equal treatment of coordinates and cells like in Array Algebra is possible.

Li and Wang [10] formalize a multidimensional model for OLAP. Core is an algebraic query language called grouping algebra which treats arrays as sets of relations plus an associated cell value which must be scalar. Operations on arrays are *add dimension*, *transfer*, *union*, *aggregation*, *rc-join* (relation/cube join), and *construct* (build array from relation). The algebra includes relations so that it can be seen as an extension of relational algebra. The model is very powerful, particularly in grouping, ordering, and aggregation. In [11] an algebra and calculus for multidimensional OLAP is presented. A multidimensional tabular database is defined as a set of tables. The model is close to the way OLAP arrays are mapped to relational star schemata. No direct mechanism is provided for join and aggregation; as by definition all first-order definable classification and aggregation functions are incorporated, these constructs can be expressed, too. Implementation of the model relies on an SQL mapping. Further important recent work in the field is described in [26, 27, 28].

In [12], an array query language, AQL, relying on Lambda calculus is presented which is geared towards scientific computing. AQL offers powerful operations on multidimensional arrays, with only slightly less generality in the aggregation mechanism than Array Algebra. The model has been implemented as a front end for querying arrays maintained in files using a geo scientific data exchange format.

An Array Manipulation Language, AML, is introduced in [13]. Two operators serve to subsample and interleave, resp., arrays based on bit sequences governing cell selection. The third operator, APPLY, corresponds to induce operations modulo the bit pattern for cell selection. Bit patterns are modeled in Array Algebra through 1-D bit arrays executing the same control function. AML is more restricted than Array Algebra in that such control arrays cannot contain arbitrary values (e.g., weights), and moreover are constrained to 1-D. According to the authors, the main application area of AML is seen in imaging.

In summary, most array frameworks nowadays are geared towards OLAP tasks, without regarding, e.g., spatio-temporal array application fields. Conversely, frameworks such as AML aiming at imaging do not consider OLAP. Sometimes array iteration or aggregation retracts to user-defined functions which, in an implementation, makes optimization difficult. All operations such as aggregation and spatial join found in these approaches are expressible in Array Algebra, too, except that dimension hierarchies usually are supported by convenient mechanisms, a feature still to be included in Array Algebra.

## 5 Conclusion

Let us be philosophic for a moment. Loosely speaking, for every data abstraction a corresponding basic operation exists. For instance, component access by name corresponds with records (in C/C++: “structs”), whereas linear traversal relates to lists and sets. For arrays, usually (nested) loops are seen as the “natural” operation, exploiting the per dimension linear ordering of cell indices. Ordering, however, is not the essence. Instead, the neighborhood defined by the indices is the crucial property: consolidation of a data warehouse cube, say, to derive weekly figures from daily data, involves seven neighbored array cells for every derived cell value. Likewise, edge detection in a 2-D raster image involves an  $n \times n$  neighborhood of each pixel for computation of the result pixels. Hence, we claim that not iteration is the operation characteristic for arrays, but (conceptually) simultaneous computation of all result array cells, in general based on the evaluation of some neighborhood for each cell.

RasDaMan Array Algebra has been designed as an algebraic framework for multi-dimensional arrays of arbitrary dimension and base type. Essentially two functionals and a sorter are sufficient for a broad range of statistics, OLAP, and imaging operations. They are declarative by nature and do not prescribe any iteration sequence, thereby opening up a wide field for query optimization and parallelization. Array Algebra is minimal in the sense that no subset of its operations exhibits the same expressive power. It is also closed in application: any expression is either of a scalar or an array type. Finally, Array Algebra does not rely on any external array handling functionality (“user-defined functions”) aside of the operations coming with the algebraic structure of the cell type. By making all operations explicit, query optimization is eased considerably.

Array Algebra comprises the formal basis for the domain-independent array DBMS RasDaMan [17] vendored by Active Knowledge GmbH<sup>3</sup>. The query language RasQL supports declarative array expressions embedded in standard SQL-92. The array query optimizer relies on about 150 algebraic rewrite rules on logical and physical level [21]. Streamlined storage management allows to distribute arrays across heterogeneous storage media. RasDaMan is being used, e.g., in the European Human Brain Database<sup>4</sup> for WWW-based access to 3-D human brain images.

Future work on the conceptual level will encompass domain-specific features such as dimension hierarchies for OLAP, including symbolic dimension handling instead of the pure numbering scheme used now, and vector/raster integration for geo applications. On the architectural level, extending RasQL functionality based on further benchmarking [17] will keep us busy for some time.

---

<sup>3</sup> See [www.active-knowledge.de/](http://www.active-knowledge.de/)

<sup>4</sup> see [www.dhbr.neuro.ki.se/ECHBD/Database/](http://www.dhbr.neuro.ki.se/ECHBD/Database/)



## Acknowledgements

RasDaFolks – Andreas Dehmel, Paula Furtado, Roland Ritsch, and Norbert Widmann – form a group with outstanding ingenuity and team spirit; it needs such people to bring a system like RasDaMan into existence.

The RasDaMan project has been partially sponsored by the European Commission under grant no. 20073.

## References

1. R. Kimball: *The Data Warehouse Toolkit*. John Wiley & Sons, 1996
2. R. Haskin, R. Lorie: *On Extending the Functions of a Relational Database System*. Proc. ACM SIGMOD, Orlando, USA, June 1982, pp. 207 - 212.
3. R. Agrawal, A. Gupta, S. Sarawagi: *Modeling Multidimensional Databases*. Proc. 13<sup>th</sup> ICDE, Birmingham, UK, April 1997, pp. 232 – 243.
4. J. Gray, A. Bosworth, A. Layman, H. Pirahesh: *Data Cube: A Relational Aggregation Operator Generalizing Group-By, Cross-Tabs, and Sub-Totals*. Technical Report MSR-TR-95-22, Microsoft Research, Advance Technology Division, Microsoft Corp., November 1995.
5. A. Bauer, W. Lehner: *The Cube-Query-Language for Multidimensional Statistical and Scientific Database Systems*. Proc. 5<sup>th</sup> DASFAA 1997, Melbourne, Australia, April 1997, pp. 263 - 272.
6. A. Shoshani: *Statistical Databases: Characteristics, Problems, and some Solutions*. Proc. VLDB 8, Mexico City, Mexico, September 1982, pp. 208 - 222.
7. A. Shoshani: *OLAP and Statistical Databases: Similarities and Differences*. Proc. PODS 16, Tucson, USA, May 1997, pp. 185 - 196.
8. E.F. Codd, S.B. Codd, C.T. Salley: *Providing Olap (On-Line Analytical Processing) to User-Analysts: An IT Mandate*. White paper, 1995, URL: [www.arborsoft.com/papers/coddTOC.html](http://www.arborsoft.com/papers/coddTOC.html)
9. L. Cabibbo, R. Torlone: *A Logical Approach to Multidimensional Databases*. EDBT 1998.
10. C. Li and X.S. Wang: *A data model for supporting on-line analytical processing*. Proc. CIKM, Rockville, USA, November 1996, pp. 81 - 88.
11. Marc Gyssens, Laks, V.S. Lakshmanan: *A Foundation for Multi-Dimensional Databases*. Proc. VLDB, Athens, Greece, August 1997, pp. 106 - 115.
12. L. Libkin, R. Machlin, L. Wong: *A Query Language for Multidimensional Arrays: Design, Implementation, and Optimization Techniques*. Proc. ACM SIGMOD, Montreal, Canada, June 1996, pp. 228 – 239.
13. A. P. Marathe, K. Salem: *A Language for Manipulating Arrays*. Proc. VLDB, Athens, Greece, August 1997, pp. 46 – 55.
14. P. Buneman: *The Discrete Fourier Transform as a Database Query*. Technical Report MS-CIS-93-37, University of Pennsylvania, 1993.
15. P. Baumann: *On the Management of Multidimensional Discrete Data*. VLDB Journal 4(3)1994, Special Issue on Spatial Database Systems, pp. 401 – 444.
16. P. Baumann: *Language Support for Raster Image Manipulation in Databases*. Proc. Int. Workshop on Graphics Modeling and Visualization in Science & Technology, Darmstadt, Germany, April 1992.
17. N. Widmann, P. Baumann: *Performance Evaluation of Multidimensional Array Storage*

Techniques in Databases. Proc. IDEAS, Montreal, Canada, June 1999 (accepted for publication)

18. G. Ritter, J. Wilson, J. Davidson: *Image Algebra: An Overview*. Computer Vision, Graphics, and Image Processing, 49(1)1994, pp. 297-336.
19. M. Stonebraker, D. Moore: *Object-Relational DBMSs: The Next Great Wave*. Morgan Kaufmann Publishers, 1996.
20. R.G.G. Cattell: *The Object Database Standard: ODMG-93*. Morgan Kaufmann Publishers, 1996.
21. R. Ritsch, P. Baumann: *Optimization and Evaluation of Array Queries*. RasDaMan Project Technical Report, FORWISS 1998.
22. M. Blaschka, C. Sapia, G. Höfling, B. Dinter: *Finding Your Way through Multidimensional Data Models*. Proc. Workshop on Data Warehouse Design and OLAP Technology DWDOT, Vienna, Austria, August 1998.
23. K.E. Iverson: *A Programming Language*. John Wiley & Sons, Inc., New York, 1962
24. J.A. Brown, H.P. Crowder: *APL2: Getting Started*. IBM Systems Journal, Vol. 30, No. 4, 1991, pp. 433 - 445
25. S. Vandenberg, D. DeWitt: Algebraic Support for Complex Objects with Arrays, Identity, and Inheritance. Proc. ACM SIGMOD, Denver, USA, May 1991, pp. 158 – 167.
26. A. Datta, H. Thomas: A Conceptual Model and an Algebra for Online-Analytical Processing in Data Warehouses. Proc. WITS 1997
27. W. Lehner: *Modeling Large Scale OLAP Scenarios*. Proc. 6<sup>th</sup> EDBT, Valencia, Spain, March 1998, pp. 153 - 167
28. P. Vassiliadis: *Modeling Multidimensional Databases, Cubes, and Cube Operations*. Proc. 10<sup>th</sup> SSDBM, Capri, Italy, July 1998