

---

1  
2  
3  
4 **SODA: A Framework for Spatial Observation Data**  
5 **Analysis**  
6

7  
8 **Sebastián Villarroya · José R.R.**  
9 **Viqueira · Manuel A. Regueiro · José**  
10 **A. Taboada · José M. Cotos**  
11

12  
13  
14  
15 Received: date / Accepted: date  
16

17  
18 **Abstract** Very large amounts of geospatial data are daily generated by many  
19 observation processes in different application domains. The amount of pro-  
20 duced data is increasing due to the advances in the use of modern automatic  
21 sensing devices and also in the facilities available to promote crowdsourcing  
22 data collection initiatives. Spatial observation data includes both data of con-  
23 ventional entities and also samplings over multi-dimensional spaces. Existing  
24 observation data management solutions lack declarative specification of spatio-  
25 temporal analytics. On the other hand, current data management technologies  
26 miss observation data semantics and fail to integrate the management of en-  
27 tities and samplings in a single data modeling solution. The present paper  
28 presents the design of a framework that enables spatio-temporal declarative  
29 analysis over large warehouses of observation data. It integrates the manage-  
30 ment of entities and samplings within a simple data model based on the well  
31 known mathematical concept of function. Observation data semantics are in-  
32 corporated into the model with appropriate metadata structures.  
33  
34

---

35 Sebastián Villarroya · José R.R. Viqueira · Manuel A. Regueiro · José A. Taboada · José M.  
36 Cotos

37 Computer Graphics and Data Engineering Group (COGRADE)  
38 Centro Singular de Investigación en Tecnoloxías da Información (CITIUS)  
39 Universidade de Santiago de Compostela  
40 Santiago de Compostela, Spain  
41 E-mail: sebastian.villarroya@usc.es

42 José R.R. Viqueira  
43 E-mail: jrr.viqueira@usc.es

44 Manuel A. Regueiro  
45 E-mail: manuelantonio.regueiro@usc.es

46 José A. Taboada  
47 E-mail: joseangel.taboada@usc.es

48 José M. Cotos  
49 E-mail: manel.cotos@usc.es  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

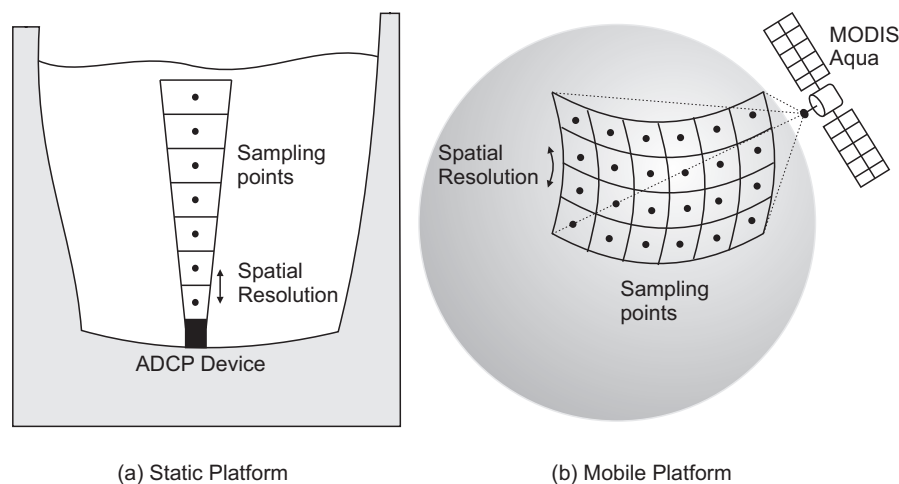
**Keywords** Spatial Data, Observation Data, Sensor Data, Data Analysis, Data Warehouse

## 1 Introduction

Nowadays, an increasing number of automatic data acquisition devices (sensors) are observing every day more variables in every day more applications domains (home automation, industrial process monitoring, health care, environmental monitoring, etc.). In many cases, the location of each observation in some reference space is an important piece of metadata that must be used during data analysis. This is always the case in applications in the area of environmental data management. As an example, a huge amount of data is generated on a daily basis by earth observation sensors on board of different satellites. Together with social media text data, the above observation data is one of the main data sources subject for current and future application of consolidated and emerging Big Data technologies.

According to the Observations and Measurements (O&M) conceptual schema [15], properties of *Entities* are either exact values assigned by some authority (name of a river, geometry of a municipality, etc.) or estimated by some (Observation) *Process* (temperature of sea surface, weight of a person, etc.). In order to adequately interpret values of *Observed Properties*, relevant observation metadata has to be recorded. Thus, it is usually mandatory to register characteristics of the specific *Process* used to generate the value. It is also mandatory to know the time instant at which the observed value starts to be valid for the *Entity* (*phenomenonTime* in [15]). Notice for example that the analysis of water (*Process*) from a river (*Entity*) might be performed in a laboratory some days after the water was obtained (*phenomenonTime*). The above metadata provides observation semantics to such property values. Observation *Processes* of very different nature might be found in real application scenarios, including physical devices (sensors), tasks performed by people and data processing algorithms. A classification of *Processes* is proposed in [32]. Two specific characteristics of a *Process* determine the type of observation data that it produces.

- A *Time-triggered Process* is performed at some predefined time frequency and therefore the observation data that it produces has the form of a regular sampling in the temporal domain. An example of this is the sampling of air temperature obtained by the temperature sensor of a meteorological station every ten minutes. On the other hand, *Event-triggered Processes* might start at any time instant, being fired by some event. For example, temperature and viscosity of volcanic lava might be measured by a volcanologist at any moment.
- If we restrict to sensors (physical *Processes* according to [32]), in-situ sensors observe *Entities* that are located in their same spatial position and they generate a single observed value at each time instant. An example of an in-situ sensor is a temperature sensor that might be installed in a



**Fig. 1** Illustration of 1D and 2D Spatial Samplings

meteorological station (static platform) or in a radio-sounding device (mobile platform). On the other hand, remote sensors observe *Entities* that are far away from its location and they use to generate various observed values (one for each observed *Entity*) at each time instant. An example of a remote sensor installed in a static platform is an Acoustic Doppler Current Profiler (ADCP), that produces at each time instant a 1D sampling of water current velocities along consecutive discrete locations of a straight line profile, either horizontal or vertical. Fig. 1(a) illustrates a vertical section of a body of water in whose bottom a ADCP sensor is installed. An example of a remote sensor installed in a mobile platform is the Moderate-resolution Imaging Spectroradiometer (MODIS) sensor installed in Terra and Aqua NASA satellites (see Fig. 1(b) for an illustration). The produced data includes a 2D regular sampling of sea surface temperatures (with a spatial resolution of 4 km) every 8 days. 2D regular samplings are called *Rasters* in the area of geographic data management.

As a consequence of the above, to effectively manage observation data, a system must provide the following general functionalities.

1. Support for the management of conventional Entity/Relationship (ER) data related to non-observed properties of *Entities*.
2. Support for observation data semantics provided by relevant observation metadata of *Observed Properties* of *Entities*. Thus, *Entity Types* together with their relevant conventional and *Observed Properties* and *Process Types* of any type should be declared to the system with the help of some declarative observation data definition language. Observation metadata of a given observed value should be provided during data insertion, including the instance of the *Process Type* (for example a specific sensing device) used to get the value and its phenomenonTime.

3. Support for the management of sampled data over temporal, spatial (1D and 2D) and spatio-temporal domains.

To the best of these authors knowledge, none of the available technologies and approaches found in data management literature provide support for all the above functionalities. In particular, observation data semantics are only explicitly supported by standards of the Open Geospatial Consortium (OGC), Sensor Web Enablement (SWE) initiative [11,15,32] and by specific observation data models and ontologies [10,29,14]. However, the support of declarative analysis over observation data is out of the scope of all those models and standards. Currently available Geographic Information System (GIS) tools [30] provide support for the recording and processing of conventional and spatial data, including Rasters, however, they lack support for declarative data analysis. Various systems have been developed for the declarative management of data streams of sensor data [28,18], however, sampled data is not supported in these systems. Many research approaches have been proposed in the area of spatial databases [20,27] and relevant functionality has been added to ISO SQL standard [24], which is implemented by well known DBMSs [31]. Currently, these tools provide support for declarative querying of spatial data, including limited support for 2D Rasters. Spatial extensions have also been implemented in NoSQL [2] and high performance Data Warehouse [23] tools. However, sampled data is not supported in these systems. Declarative analysis over very large collections of sampled data is supported by array data managers [8,12]. However, declarative analysis of relational data is not user friendly with array data structures. Finally, the integrated management of relational and array data is attempted in [39]. To achieve this, the user has to deal with both relational and array semantics, which in these authors opinion is not user friendly.

Based on all the above limitations of currently available systems and approaches, the objective of the present work is the design of SODA (Spatial Observation Data Analysis), a framework for declarative spatio-temporal analysis in very large warehouses of spatial observation data. The framework seamlessly integrates entity based and sampling data in a simple data model based on the well known mathematical concept of function. The model incorporates also observation data semantics. A main handicap is that it is not based on relational or object oriented paradigms, however the data model is very simple and the language combines logical and functional constructors already present in other well known languages like XQuery. The contributions of this paper may be summarized as follows.

- Formalization of a data model for the integrated management of entity and sampling data, using a new hybrid logical-functional paradigm.
- Definition of a spatio-temporal declarative data analysis language for the above data model.
- Definition of a data warehouse data model with support for observation data semantics. Application of the above language to the declarative def-

1       initiation of new observation *Processes* that are executed by the framework  
2       during observation data load.

- 3  
4       – Brief discussion of physical level issues related to the column-oriented im-  
5       plementation of the framework that is currently being undertaken. The im-  
6       plementation must exploit parallelization in current multi-core hardware  
7       architectures.

8       The remainder of this paper is organized as follows. Section 2 discusses and  
9       evaluates technologies and research approaches related to the present one. The  
10       formalization of the data model for observation data warehouses is given in  
11       Section 3. Section 4 describes the spatio-temporal analysis language provided  
12       by the framework. Column-oriented implementation issues are discussed in  
13       Section 5. Finally, Section 6 concludes the paper and outlines pieces of further  
14       work.  
15

## 16 17 18 **2 Related Work**

19  
20       Various related approaches and technologies are now described and compared  
21       with respect to the following criteria, which is derived from the generic func-  
22       tionalities that an observation management system must support.

- 23  
24       1. *Direct support of observation semantics*: In order to perform effective anal-  
25       ysis and interpretation of values of *Observed Properties* of *Entities*, some  
26       important metadata have to be recorded and linked to the observed values.  
27       In particular, at least, it is important to record a reference to the specific  
28       *Process* used to generate the observation and also the time at which the  
29       observed value applies to the *Entity*. *Process* instances have to be clas-  
30       sified into *Process Types* as *Entities* are also classified into *Entity Types*  
31       in the classical E/R model. Besides, the system should also support the  
32       recording of properties of *Processes*. Thus, for example, in a meteorologi-  
33       cal observation domain, a *Process Type* “TemperatureSensor” could have  
34       self described properties “DeviceId” and “InstallationDate”. Each value of  
35       the *Observed Property* “AirTemperature” of each Meteorological Station  
36       *Entity* should be linked to the specific instance of “TemperatureSensor”  
37       that was used to measure it.  
38       2. *Support for the management of sampled data*: Beyond the classical E/R  
39       data, an observation data management approach must also provide data  
40       structures and operations that enable the efficient processing of sampled  
41       data. As it was already reported in the introduction, temporal samplings  
42       are generated by *Time-triggered Processes*, whereas spatial samplings are  
43       usually produced by remote sensors. It is noticed that the use of classical  
44       relational-based models for sampled data results in either highly inefficient  
45       approaches or complex nested models as will be shown below in this section.  
46       3. *Support for multi-resolution temporal and spatial data*: The observation  
47       data generated by currently available sensors is produced with different  
48       temporal and spatial resolutions. Thus for example, MODIS generates sea  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 surface temperature data with a temporal resolution of 8 days and with  
2 a spatial resolution of 4 km. An observation data management approach  
3 should provide a data type system that simplifies the transformation be-  
4 tween different temporal and spatial resolutions during the evaluation of  
5 operations.

- 6  
7 4. *Simple data modeling approach*: In the context of this evaluation, a sim-  
8 ple data model is the one that uses just one non-nested data structure.  
9 Thus, for example, a relational model is considered simple, whereas an  
10 object-relational one is not, as nested arrays and collections are supported.  
11 Besides, data models that use different data structures with different se-  
12 mantics for different types of data are also considered non-simple. It is  
13 obvious that the efficient implementation of a nested data model is far  
14 more complicated than the implementation of a non-nested one. On the  
15 other hand, it is also clear that having to deal with various data structures  
16 with different semantics leads to interfaces that are not user friendly.
  - 17 5. *Model based on a well known paradigm*: The definition of data models  
18 that are based on well known paradigms as the relational one allows to  
19 take advantage of many years of user experience, improving their learning  
20 curve.
  - 21 6. *Stream Processing approach*: If there are real-time requirements and the  
22 amount of data to be recorded is not large then the approach must inte-  
23 grate the efficient processing of input data streams with small stored data  
24 structures to produce output data streams. Stream processing approaches  
25 are commonly known as Complex Event Processing (CEP) (Information  
26 Flow Processing Systems in [16]) and they rely on the evaluation of Con-  
27 tinuous Query Language (CQL) expressions [7,25].
  - 28 7. *On Line Transaction Processing (OLTP) approach*: If real-time require-  
29 ments are present with simpler temporal patterns but large amounts of data  
30 have to be recorded, then an OLTP processing approach is required. This  
31 approach is traditionally supported by conventional DBMSs for reasonably  
32 large data collections and provided by both NoSQL [2,1] and NewSQL [4]  
33 solutions in the new era of Big Data Management.
  - 34 8. *On Line Analytical Processing (OLAP) approach*: Non-real time analyt-  
35 ics over very large data sets is supported by OLAP approaches provided  
36 by Business Intelligence solutions over Data Warehouse technologies. High  
37 performance implementations include Hewlett-Packard Vertica [3], which  
38 is an evolution of CStore [36] and the open source Monetdb database [23].  
39 The efficient implementations of these Big Data solutions are based on  
40 recent research on column-oriented technologies. The key feature of these  
41 approaches is that relational data is recorded by columns, instead of the  
42 classical row storage. This enables on the one hand the application of ef-  
43 ficient compression techniques to the data and even to perform some pro-  
44 cessing over compressed data and on the other hand avoids retrieving from  
45 storage columns that are not involved in computations. The main draw-  
46 back is that insertions, updates and deletions of data are not efficient, that  
47 is why they are suitable for data warehouses.
- 48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1 9. *Support for declarative processing*: In data management, the advantages of  
2 a declarative language like SQL over a procedural approach are very well  
3 known. This is a clear motivation for trying to apply data management  
4 technologies in some application domains where procedural solutions are  
5 dominant. The management of sampled observation data in environmental  
6 applications is one such domain.
- 7 10. *Support for aggregation*: Aggregation functionalities through statistical meth-  
8 ods are at the kernel of OLAP and must be supported to effectively perform  
9 observation data analysis.
- 10 11. *Support for iterative processing*: Recursive queries are required in only few  
11 data management applications. This is the reason why such functionalities  
12 were out of the scope of first SQL implementations. Current ISO SQL stan-  
13 dard and DBMSs vendors support a kind of limited recursion. Regarding  
14 the analysis of observation data in environmental applications, such func-  
15 tionalities are commonly required to perform many simulations. Examples  
16 of these are forest fire propagation, oil spills, flooding, etc. Thus, although  
17 it is not a kernel functionality, the support for iterative processing is a  
18 desirable feature.
- 19 12. *Data processing based on a well known language*: As in the case of the data  
20 model, the definition of query languages that are based on well known ones  
21 as SQL is a clear advantage.
- 22 13. *Availability of efficient implementation*: A data management approach is  
23 really useful if it can be efficiently implemented. A prototype implementa-  
24 tion demonstrates the viability of the approach and its use in real applica-  
25 tion domains shows its matureness.

26  
27  
28  
29 Based on the above evaluation criteria various research approaches and  
30 available technologies are now classified and qualitatively compared, including  
31 also the present SODA framework. An overview of such comparison is given  
32 in the table of Fig. 2, where each approach is marked with “Y” or “P” if  
33 it, respectively, supports or partially supports the relevant criterion. A more  
34 detailed discussion is given below.

35  
36 *OGC SWE Standards*: The Sensor Web Enablement (SWE) of the Open  
37 Geospatial Consortium (OGC) provides a series of standards for the interfaces  
38 of web services related to the management of environmental observation data.  
39 In particular, the Observations and Measurements (O&M) [15] and Sensor  
40 Model Language (SensorML) [32] were already mentioned in the introduction.  
41 The Sensor Observation Service (SOS) [11] defines a web service interface to  
42 query observation data collections, either stored or directly obtained from the  
43 devices. Data is transferred between client and server in standard XML encod-  
44 ings of O&M and SensorML models. Query capabilities of SOS are limited  
45 to just filtering. Regarding data processing, OGC defines the Web Processing  
46 Service (WPS) [34] interface that enables the invocation of data processing  
47 algorithms through the web. Various implementations of the above standards  
48 exist already in the market, both with commercial and open source licenses. In  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

	Obs. Semantics	Sampled Data	Multiresolution	Simple Model	Well Known Model	Stream Proc.	OLTP	OLAP	Declarative Proc.	Aggregation	Iterative Proc.	Well Known Lang.	Impl. Available
OGC SWE Stds	Y	Y			Y								Y
Obs. Data Models	Y			Y	Y								
GIS		Y		Y	Y								Y
Sensor Stream				Y	Y	Y			P	P		P	Y
Spat. and ST DBMSs		Y			Y		Y	Y	P	P	P	P	Y
Spatial NoSQL							Y		Y				Y
Spatial HP DW				Y	Y			Y	P	P		P	Y
Array Data Managers		Y		Y				Y	Y	Y			Y
SciQL		Y			P			Y	Y	Y		P	Y
SODA	Y	Y	Y	Y				Y	Y	Y			

**Fig. 2** Comparison of related approaches and technologies

general it is obvious that O&M provide appropriate support for the modeling of observation semantics and sampled data. Different spatial and temporal resolutions are supported but transformations are a user matter. The underlying object oriented data modeling approach with XML encodings is well known, however, to support sampled data nested structures are required. Declarative data processing is not supported at all as WPS just provides means for remote procedure calls.

*Observation Data Models:* Beyond the above O&M OGC standard, various data models and ontologies have been proposed to support observation data semantics [10,29,14]. They are based on well known paradigms and provide observation data semantics with simple data modeling approaches. However, sampled data and multi-resolution is out of the scope of these models as it is also any kind of data processing.

*Geographic Information Systems (GIS) :* Currently, a wide variety of GIS tools, both with commercial and open source licenses, are available. A representative example of them is GRASS [30], which supports the management of any kind of geographic data, including Rasters, recorded in many different well known models and formats. Raster data management is usually formalized with relevant Raster algebras [13]. Observation semantics are not considered in

1 GIS and although the managed data may have many different spatial resolu-  
2 tions, transformations between them have to be explicitly done by the user to  
3 perform operations. Spatial data processing is a strength of tools like GRASS,  
4 however, it is performed by the execution of a very large amount of different  
5 commands, therefore, a declarative language is missing. Notice that the user  
6 must know which is the functionality of each command and how to combine  
7 them, thus only real experts may take real advantage of spatial data analysis  
8 with GIS tools.  
9

10  
11 *Sensor Stream processing approaches:* Various Stream Processing approaches  
12 have been explicitly proposed for the management of data generated by sensor  
13 networks [28,18]. Despite of being defined for sensor data, observation data  
14 semantics are not explicitly incorporated by the system and are delegated  
15 to user interpretation. Any kind of spatial data management is out of the  
16 scope of these approaches. They support declarative real-time processing of  
17 streams with aggregation functionality based on SQL like languages. Real-  
18 time requirements of these approaches are clearly in conflict with the support  
19 of iterative processing.  
20  
21

22  
23 *Spatial and Spatio-temporal DBMSs:* Many temporal extensions have been  
24 proposed for the classical relational model [17,35]. Recently, some characteris-  
25 tics have been incorporated into ISO SQL standard [26]. Various spatial [20,27,  
26 37] and spatio-temporal [21,38] extensions to classical models have been pro-  
27 posed in the literature. Spatial functionality has already been added to ISO  
28 SQL standard [24], which is currently implemented by most of the available  
29 DBMSs (see [31] for an example). The direct support of observation semantics  
30 is out of the scope of spatial DBMSs. They support the management of conven-  
31 tional E/R data with a well known object-relational paradigm and SQL, where  
32 properties of entities might have spatial data types (point, line, surface, etc.).  
33 Extensions for Raster data are also supported by some approaches and sys-  
34 tems [31], however, they require nested structures and do not provide explicit  
35 support for multi-resolution. They can be used both for OLTP and OLAP,  
36 but in the general case they were not designed with Big Data requirements in  
37 mind. Regarding declarative data processing, it is only efficiently supported  
38 for non-sampled data and it includes both aggregations and SQL recursion  
39 for iterative queries. To manipulate raster data with SQL constructors it has  
40 to be unnested from a complex value of a raster data type, which is a highly  
41 inefficient task.  
42  
43

44  
45 *Spatial NoSQL approaches:* Systems following a NoSQL approach and pro-  
46 viding spatial data management capabilities are still few. An example is the  
47 extension of MongoDB [2] with support for the management of GeoJSON en-  
48 coded data. Their functionality is very limited both in data modeling and  
49 processing.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 *Spatial High Performance Data Warehouse approaches:* To the best of these  
2 authors knowledge, only the Monetdb DBMS [23] provides spatial functional-  
3 ity on top of a high performance column-oriented implementation for OLAP.  
4 Neither observation semantics nor sampled data are directly supported by the  
5 system. Therefore, declarative processing is only partially supported. It lacks  
6 recursive queries, therefore iterative processing is not supported.  
7

8  
9  
10 *Array Data Managers:* The management of sampled data fits very well ar-  
11 ray data management approaches. Currently, array algebras like the one in  
12 [9] are the basis for the development of relevant array data managers [8,12].  
13 These systems use a simple array data model to provided high performance  
14 OLAP over very large arrays. They provide declarative array query languages  
15 that include aggregation capabilities. Iterative processing is not supported.  
16 Although the flavor of the languages is similar to that of SQL, both the ar-  
17 ray semantics and complex array operators make them quite cumbersome for  
18 DBMS users. Observation semantics are out of the scope of these systems and  
19 multi-resolution is not explicitly supported.  
20

21  
22  
23 *SciQL:* A SQL-based query language for science applications SciQL is defined  
24 in [39]. The language enables the integrated analysis of relational and ar-  
25 ray data, therefore both *Entities* and sampled data are supported. Obviously,  
26 declarative query with aggregation is supported, however, current implemen-  
27 tation with MonetDB [23] technologies does not include recursion for iterative  
28 processing support. The incorporation of an array data structure adds com-  
29 plexity to the relational model and the extension of SQL with array semantics  
30 makes it less friendly to DBMS users. Observation semantics are out of the  
31 scope of the approach and specific support for multi-resolution in the type  
32 system is not provided.  
33

34  
35  
36 *SODA:* The SODA framework described in the current paper is also added  
37 to the comparison for qualitative evaluation purposes. As it is shown, SODA  
38 supports both observation semantics and sampled data. It combines a multi-  
39 resolution type system with a simple non-nested data structure based on the  
40 well known mathematical concept of function. Declarative spatio-temporal  
41 analysis is supported and aggregation functionality is also incorporated. Sup-  
42 port for iterative processing functionality and the efficient implementation of  
43 SODA are part of future work. Both the data model and query language of  
44 SODA are not based on the well known relational paradigm, which is somehow  
45 a setback for current DBMS users. However, the formalism is simple and it is  
46 based on the well known mathematical concept of function. The declarative  
47 language combines logical and functional constructors very similar to those al-  
48 ready present in languages like XQuery. These characteristics will be described  
49 throughout the paper.  
50

51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

### 3 Observation Data Warehouse

This section describes the observation data warehouse that supports observation data recording. The general functionalities of an observation data management system briefly described in the introduction give rise to the following requirements for an observation data warehouse.

1. The model must support the representation of classical E/R data integrated with temporal, spatial and spatio-temporal sampled data.
2. The model must provide direct support for observation semantics, through the representation of appropriate required metadata, including *Processes* and *phenomenonTime*.
3. The type system of the model must provide support for the representation of temporal and spatial data at different resolutions. Besides, the transformation between those resolutions must also be simplified with relevant implicit and explicit type castings.

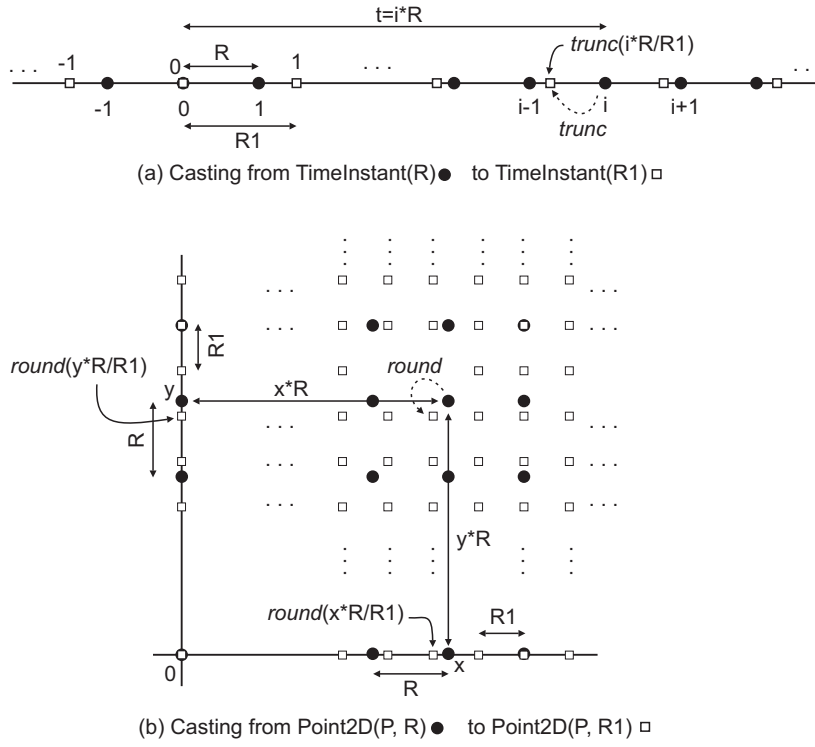
Based on the above requirements, an underlying spatio-temporal data model is first defined and on top of it structures for observation metadata are added to achieve the final data model for the observation data warehouse.

#### 3.1 Underlying Spatio-temporal Data Model

As it was already shown in the previous section, supporting sampled data with a relational formalism leads to highly inefficient approaches. On the other hand, the use of functional models for the management of E/R data has already been tried in the area of Functional Databases [19]. The proposed spatio-temporal data model is based on the well known mathematical concept of function. It consists of conventional, temporal and spatial data types, functions to manipulate data values called *Intensional Mappings* and functions to record data values called *Extensional Mappings*. The recording of data singletons, called *Constants* is also supported by the model.

##### 3.1.1 Data types

Conventional data types consist of *Boolean*, *CString* (variable size character strings), *Integer* and two floating point numeric types, *Float* and *Double*. Parametric type *FixedPrecision(P, S)* enables the user to specify precision (P, maximum number of decimal digits used) and Scale (S, number of digits in the fractional part) in a fixed point numeric representation. Default and also maximum values for P and S are defined by the system. We shall denote those values as DP (default P), MP (maximum P), DS (default S) and MS (maximum S). Any value N of type *FixedPrecision(P, S)* can be written in the form  $n * 10^{-S}$ , where  $n$  is an integer in the range  $(-10^P, 10^P)$ . All the above data types include a special *Undefined* value denoted by  $\perp$ .



**Fig. 3** Type Castings between Temporal and Spatial Types

*Temporal Data Types:* The following three data types enable the modeling of the discrete multi-resolution time values.

**Definition 1** Let  $R$  (resolution) be a value of data type  $\text{FixedPrecision}(P, S)$ . Then the following temporal data types are defined.

$\text{TimeInstant}(R)$ :

$$\{t * R \mid t \in \text{Integer} \wedge -10^{MP} < t < 10^{MP}\} \cup \{\perp\}$$

$\text{Time}(R)$ :

$$\{t * R \mid t \in \text{Integer} \wedge 0 \leq t * R < 24\text{hours} * 3600\text{seconds/hour}\} \cup \{\perp\}$$

$\text{Date}(R)$ :

$$\text{TimeInstant}(86400)$$

Time instant values are semantically interpreted as temporal shifts (positive or negative) in seconds from a reference system defined time instant ( $t = 0$ ). A common value for such a reference time instant in current DBMS implementations is “1970-01-01 T 00:00:00.000000”.

*Spatial Data Types:* The following two data types enable the modeling of 1D and 2D Euclidean spaces.

**Definition 2** Let  $P$  (Precision) and  $R$  (Resolution) be values of types *Integer* and *FixedPrecision*(PR, SR), respectively. Then the following spatial data types are defined.

*Point1D*(P, R):

$$\{x * R \mid x \in \text{Integer} \wedge -10^P < x < 10^P\} \cup \{\perp\}$$

*Point2D*(P, R):

$$\{(x * R, y * R) \mid x, y \in \text{Integer} \wedge -10^P < x, y < 10^P\} \cup \{\perp\}$$

*Geometric Data Types:* Based on data type *Point2D*(P, R) and on the standard specification given by [24], the following data types enable the modelling of geometries in 2D Euclidean spaces.

- *LineString*(P, R): Vector polylines defined by sequences of elements of *Point2D*(P, R).
- *Polygon*(P, R): Vector polygons, possibly with holes, whose borders are defined by sequences of elements of *Point2D*(P, R).
- *GeometryCollection*(P, R): Heterogeneous collections of geometries of any of the data types *Point*(P, R), *Polyline*(P,R) and *Polygon*(P,R).
- *MultiPoint*(P, R): Homogeneous collections of geometries of type *Point2D*(P,R).
- *MultiLineString*(P, R): Homogeneous collections of geometries of type *LineString*(P,R).
- *MultiPolygon*(P, R): Homogeneous collections of geometries of type *Polygon*(P,R).
- *Geometry*(P, R): Abstract type. Geometries or collections of geometries of any of the above 2D types.

### 3.1.2 Intensional Mappings

Informally, *Intensional Mappings* are functions defined over the available data types. These functions are always defined intensionally, in a mathematical sense, by either some algorithm or expression. Notice the difference with *Extensional Mappings* defined below in Subsection 3.1.3.

**Definition 3** If  $T_1, T_2, \dots, T_n$  is a possibly empty sequence of not necessarily distinct data types and  $T$  is also a data type, then an *Intensional Mapping* with signature  $M(T_1, T_2, \dots, T_n):T$  is defined as a function  $M: T_1, T_2, \dots, T_n \rightarrow T$ .

Primitive *Intensional Mappings* are directly provided by the system in the syntactic form of operators, functions and type castings. Implicit type castings are applied between types of the same category (Boolean, CString and Numeric) during the evaluation of functions and operations.

*Temporal Intensional Mappings:* Let  $t = i * R, t_1 = i_1 * R$  be two values of type  $TimeInstant(R)$  and let  $n = i_n * 10^{-S_R}$  be a value of type  $FixedPrecision(P_n, S_R)$ , where  $S_R$  matches the scale of the  $FixedPrecision$  type of  $R = r * 10^{-S_R}$ . Then the expression  $t - t_1$  yields value  $(i - i_1) * R$  of type  $FixedPrecision(MP, S_R)$ . The expression  $t - n$  yields the  $TimeInstant(10^{-S_R})$  value  $(i * r - i_n) * 10^{-S_R}$ . The expression  $t + n$  yields the  $TimeInstant(10^{-S_R})$  value  $(i * r + i_n) * 10^{-S_R}$ . A complete ordering is straightforwardly defined for values of each temporal data type and relevant comparison operations might be applied accordingly. Temporal functions are also provided although their definition is out of the scope of this paper. Castings between temporal types enable multi-resolution temporal management. Thus, the expression “ $cast(t \text{ to } TimeInstant(R_1))$ ” yields the  $TimeInstant(R_1)$  value  $trunc(i * R/R_1)$ , as it is illustrated in Fig. 3(a). Type castings are implicitly applied between temporal types during the evaluation of functions and operations.

*Spatial Intensional Mappings:* If  $p = (x * R, y * R), p_1 = (x_1 * R, y_1 * R)$  are two values of type  $Point2D(P, R)$ , then the expressions  $p + p_1$  and  $p - p_1$  yield respectively values  $((x + x_1) * R, (y + y_1) * R)$  and  $((x - x_1) * R, (y - y_1) * R)$ , both of the same  $Point2D(P, R)$  data type. For completeness, a total ordering is defined for  $Point2D$  data types, based on some space-filling curve [33]. Such an ordering enables the application of comparison operators. If  $x, y$  are values of type  $FixedPrecision(P, S)$ , then function  $point2d(x, y)$  yields the value  $(x, y)$  of type  $Point2D(P, 10^{-S})$ . Inversely, functions  $xcoord(p)$  and  $ycoord(p)$  yields respectively the values  $x * R$  and  $y * R$  of type  $FixedPrecision(P, S_R)$ , where  $S_R$  is the scale value of the  $FixedPrecision$  type of  $R$ . Other spatial functions are provided, including *distance* and *direction*, however, their definition is out of the scope of the paper. Type castings are defined between  $Point2D$  types that enable transformations between different spatial resolutions. Thus, the expression “ $cast(p \text{ to } Point2D(P_1, R_1))$ ” yields value  $(round(x * R/R_1), round(y * R/R_1))$ , as it is illustrated in Fig. 3(b). Similar operations, functions and castings are also provided for  $Point1D$  types. Type castings are automatically applied between spatial types during the evaluation of operations and functions.

*Geometric Intensional Mappings:* Geometric functions specified in [24] are also supported. Besides, type castings enable the transformation between geometric values of different spatial resolutions.

### 3.1.3 Extensional Mappings

Informally, an *Extensional Mapping* is a function that maps values from a finite domain, defined as a finite subset of the Cartesian Product of data types, to a codomain defined by a data type. Components of the domain of an Extensional Mapping are called *Dimensions*. Informally, a *Dimension* is just a set of values of a given data type.

**Definition 4** A *Dimension*  $d$  over data type  $T$ , denoted by  $d(T)$ , is defined as a non-empty finite subset of  $T - \{\perp\}$

Finite sets of values of a given data type are recorded in *Dimensions*. Temporal and spatial samplings, defined below, are specific types of *Dimensions* of special interest for spatio-temporal data representation.

**Definition 5** Let  $min$  and  $max$ ,  $min < max$ , be two values of the same *TimeInstant*, *Time*, *Date* or *Point1D* type  $T$ . Then a *1D Sampling*  $S$  from  $min$  to  $max$ , denoted by  $S(min, max)$  is defined as the following *Dimension* over data type  $T$ .

$$\{s \in T \mid min \leq s \leq max\}.$$

**Definition 6** Let  $s_m = (x_m, y_m)$  and  $s_M = (x_M, y_M)$ ,  $x_m < x_M$  and  $y_m < y_M$ , be two values of the same *Point2D* type  $T$ . Then a *2D Sampling*  $S$  from  $s_m$  to  $s_M$ , denoted by  $S(s_m, s_M)$  is defined as the following *Dimension* over data type  $T$ .

$$\{(x, y) \in T \mid x_m \leq x \leq x_M \wedge y_m \leq y \leq y_M\}.$$

Efficient physical structures avoid having to record the values of very large samplings (see Subsection 5.2).

**Definition 7** If  $d_1, d_2, \dots, d_n$  is a sequence of not necessarily distinct *Dimensions* and  $T$  is a data type, then a *Extensional Mapping* with signature  $M(d_1, d_2, \dots, d_n): T$  is defined as a function  $M: d_1, d_2, \dots, d_n \rightarrow T$ . The definition of a *Extensional Mapping* is represented by a finite set of pairs  $(d, c)$ , where  $d$  is an element of  $d_1 \times d_2 \times \dots \times d_n$  and  $c$  is an element of  $T$ .

*Extensional Mappings* in the proposed model have an extensional mathematical definition, since each valid combination of domain and codomain values is explicitly recorded in the model. Notice the difference with *Intensional Mappings*, whose definition is not given by an exhaustive listing of elements.

Informally, a *Constant* is just one value of a given data type.

**Definition 8** A *Constant*  $C$  of type  $T$ , denoted by  $C:T$  is defined as an atomic value of type  $T$ .

Following a functional database approach [19], *Dimensions* and *Extensional Mappings* enable the modeling of Entities and Relationships between them. For example, student, courses and enrollment data may be modeled as follows.

*Dimensions*

StudentId(Integer)

CourseId(Integer)

*Extensional Mappings*

StudentName(StudentId):CString

CourseName(CourseId):CString

EnrollMentDate(StudentId, CourseId):Date

Beyond classical ER data, the model elegantly integrates the representation of scientific and sensor data. As an example, the evolution of sea surface

1 temperature and salinity data over time (sampled every minute), depth (sam-  
 2 pled every meter) and space (sampled every kilometer) might be modeled as  
 3 follows:

4 *Constants*

5 minTime:TimeInstant(600), maxTime:TimeInstant(600),  
 6 maxDepth:Point1D(1),  
 7 southWestCorner:Point2D(1000), northEastCorner:Point2D(1000)

8 *Dimensions*

9 Time(minTime, maxTime)  
 10 Depth(*cast*(0 to Point1D(1)), maxDepth)  
 11 Space(southWestCorner, northEasCorner)

12 *Extensional Mappings*

13 Temperature(Time, Depth, Space):FixedPrecision(5, 2)  
 14 Salinity(Time, Depth, Space):FixedPrecision(5, 2)

15  
 16  
 17 The above data model fits very well column-oriented DBMS implementa-  
 18 tion techniques, which are currently the trending technology for the imple-  
 19 mentation of data warehouses (see Subsections 5.2 and 5.3).

20  
 21  
 22  
 23 

### 3.2 Observation Data Model

24 Observed *Entities* and samplings are modeled in the proposed model with the  
 25 concept of *Feature*. As in the E/R model, *Features* are classified into *Feature*  
 26 *Types*. Both key and non-key *Properties* of *Features* may be defined. *Properties*  
 27 of features may be observed by *Processes*.

28 *Processes* are also classified into *Process Types* and *Properties* of them may  
 29 also be defined. For our purposes, a *External* process is performed out of the  
 30 scope of the framework and the observations that it produces are loaded into  
 31 the data warehouse with a typical ETL task. On the other hand, an *Internal*  
 32 process is executed by the framework during the ETL task to produce obser-  
 33 vations derived from both the data that are being loaded and the data already  
 34 available in the data warehouse. For our purposes, we shall classify processes  
 35 both as either *Time-triggered* or *Event-triggered* and as either *External* or  
 36 *Internal*.

37  
 38 The above concepts and the relationships between them are depicted in  
 39 the UML class diagram of Fig. 4. The schema of a spatial observation data  
 40 warehouse is defined in SODA using a XML language called XODDL (XML  
 41 Observation Data Definition Language), which is defined in the present work  
 42 as an XML encoding of the UML class diagram of Fig. 4.

43 In order to use the framework, the administrator must create a dataset  
 44 scheme using XODDL. Next, *Observation Process* metadata is added to the  
 45 framework. Internal *Processes* have also to be defined (see Subsection 4.2).  
 46 Next, Feature Types with their non-observed properties are inserted. At this  
 47 stage, the framework is ready to execute observation data ETL tasks. Each  
 48 execution appends new values stamped with appropriate temporal data (*phe-*  
 49  
 50  
 51  
 52  
 53  
 54  
 55  
 56  
 57  
 58  
 59  
 60  
 61  
 62  
 63  
 64  
 65

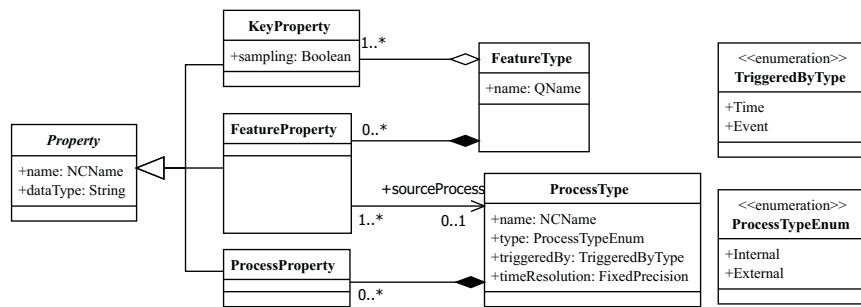


Fig. 4 UML Class Diagram of the Observation Data Model.

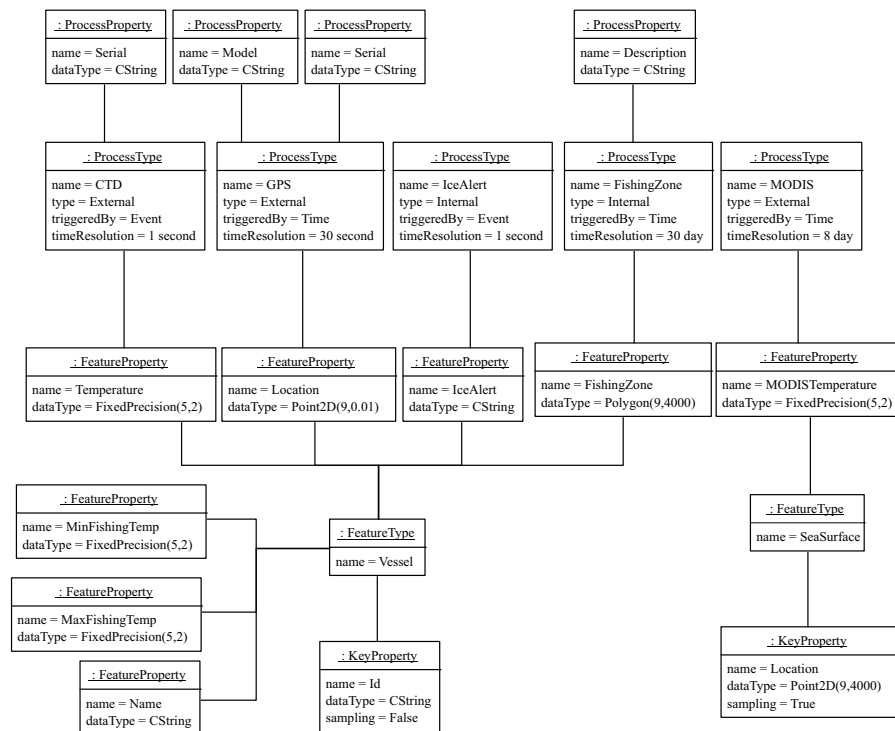


Fig. 5 UML Object Diagram of a Running Example.

*nomenonTime*) to *Observed Properties of Feature Types*. Notice that new data is always appended and deletions and updates are not supported (except for data administration purposes).

The instances of the above concepts that define a proposed running example are depicted in the UML object diagram of Fig. 5. A more detailed description of *Feature* and *Process Types* and relevant *Dimensions* and *Extensional Mappings* used to record their data is given below.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

### 3.2.1 Feature Types

*Feature Types* enable the integrated modeling of entities and samplings. Two *Feature Types* are defined in the running example, *Vessel* is used to model fishing boats (entities) and *SeaSurface* is used to model a geographic sampling. A *Feature Type* is defined by one or more *Key Properties* and zero or more *Feature Properties*. *Key Properties* are used to model either key properties of entities or dimensions of samplings. *Feature Properties* are used to model either non-key properties of entities or sampled properties. In the running example, vessels are identified by a *Vessel.Id* key property of type string, whereas the spatial dimension of *SeaSurface* is modeled by key property *SeaSurface.Loc* of type *Point2D(9, 4000)*. Notice that an attribute (*sampling="true"*) is added to the *SeaSurface.Loc* key property to indicate that it is a sampling dimension. Non-key properties of *Vessel* enable the recording of its name, minimum and maximum fishing temperatures, location, water temperature measured by the boat, ice alert risk according to measured temperature and the spatial fishing zone according to the fishing temperature interval. *SeaSurface* is a sampling that has just one *MODISTemperature* property that records the temperature at each location provided by MODIS sensor installed on board of Terra and Aqua NASA satellites. When the values of a specific *Feature Property* are generated by some *Process*, the relevant *Process Type* metadata is referenced by attribute "sourceProcessType". Notice that in the running example, only names and fishing temperatures of vessels are not generated by Observation Processes.

The following *Dimensions* and *Extensional Mappings* are recorded in the framework for each *Feature Type* *F*.

- For each *Key Property* *KP* of type *KPT* of *F*, if the attribute *sampling* is false then a non-sampling *Dimension* *F.KP(KPT)* is recorded, otherwise, a sampling *Dimension* *F.KP(lo, hi)* is recorded, where *lo* and *hi* of type *KPT* are the values that define the boundaries of the sampling recorded in the framework. In the running example, dimensions

*Vessel.Id*(CString) and  
*SeaSurface.Loc*(lo:Point2D(9,4000),hi:Point2D(9, 4000))

enable the recording of relevant key properties. Notice that in the framework, each *Dimension* and *Extensional Mapping* must have a unique name, therefore, to avoid conflicts, the name of the *Feature Type* (or *Process Type*) is added as a prefix to the name of the relevant property.

- For each *Feature Property* *FP* of type *FPT* of *F*, such that *FP* is not generated by any source *Process*, an *Extensional Mapping*

*F.FP(F.KP<sub>1</sub>, F.KP<sub>2</sub>, ... F.KP<sub>n</sub>):FPT*

is recorded, where each *F.KP<sub>i</sub>* is a *Dimension* corresponding to a relevant key property of *F*. In the running example, *Extensional Mappings*

*Vessel.Name*(*Vessel.Id*):CString,  
*Vessel.MinFishingTemp*(*Vessel.Id*):FixedPrecision(5, 2) and  
*Vessel.MaxFishingTemp*(*Vessel.Id*):FixedPrecision(5, 2)

enable the recording of relevant non-key feature properties not generated by *Observation Processes*.

- For each *Feature Property* FP of type FPT of F, such that FP is generated by an observation source process of type SP and  $KP_1, KP_2, \dots KP_n$  are key properties of F:

1. An *Extensional Mapping*

F.FP( $F.KP_1, F.KP_2, \dots F.KP_n, SP.Time$ ):FPT

enables the recording of the generated observations.

2. An *Extensional Mapping*

F.FP.Process( $F.KP_1, F.KP_2, \dots F.KP_n, SP.Time$ ):CString

enables the recording of the specific *Observation Process* used at each time instant to generate observations.

Thus, for example,

Vessel.Location(Vessel.Id, GPS.Time):Point2D(9, 0.01)

enables the recording of the evolution with respect to time of the vessel location generated by the GPS devices. The identifier of the GPS device used at each specific instant in a given vessel is provided by *Extensional Mapping*

Vessel.Location.Process(Vessel.Id, GPS.Time):CString.

### 3.2.2 Process Types

*Process Types* record metadata of referenced observation *Process*. A *Process Type* may be either *Time* or *Event-Triggered* at a given time resolution. Thus, for example, GPS processes of vessels generate vessel location observations every 30 seconds. On the other hand, at any moment the vessel crew uses a CTD sensor to obtain a temperature observation, whose time is recorded with a resolution of 1 second. The above processes are external to the framework. An example of an internal process is the IceAlert that computes an ice risk value (either “High” or “Medium”) every time the CTD sensor generates an observation with a temperature below zero. Internal processes are executed by the framework during each ETL task to generate calculated *Feature Properties*. Thus, the FishingZone property of each vessel (Polygon2D type) is generated from sea surface temperatures provided by MODIS. The capabilities of SODA to define internal analytical processes are described in Subsection 4.2. Finally, each process of a given *Process Type* might be described by a set of *Process Properties*. Thus, for example, a GPS process has Model and Serial properties as part of its metadata.

The following *Dimensions* and *Extensional Mappings* are recorded in the framework for each *Process Type* P.

- A dimension P(CString) enables the recording of all the identifiers of processes (instances) of P. In the running example, dimensions
  - GPS(CString)
  - CTD(CString)
  - IceAlert(CString)

- 1 FishingZone(CString) and  
 2 MODIS(CString)  
 3 enable the recording of all the required process identifiers.  
 4 – For each *Process Property* PP of type PPT, an *Extensional Mapping*  
 5 P.PP(P):PPT  
 6 is recorded. Thus, *Extensional Mappings*  
 7 GPS.Model(GPS):CString and  
 8 GPS.Serial(GPS): CString  
 9 enable the recording of GPS models and serial numbers in the running  
 10 example. Similarly, *Extensional Mappings*  
 11 CTD.serial(CTD):CString and  
 12 FishingZone.Description(FishingZone):CString  
 13 enable the recording of CTD serial numbers and FishingZone internal Pro-  
 14 cess descriptions.  
 15 – If P is a time-triggered Process of resolution R then a *Sampling*  
 16 P.Time(lo:TimeInstant(R), hi:TimeInstant(R))  
 17 is recorded, where lo and hi are respectively the lowest and highest time  
 18 instants inserted in the framework for observation times generated by P.  
 19 On the other hand, if P is an event-triggered Process of resolution R then  
 20 a non-sampling *Dimension*  
 21 P.Time(TimeInstant(R))  
 22 is recorded.  
 23  
 24  
 25

## 26 4 Spatial Observation Data Analysis

27  
 28 The following requirements for spatial observation data analysis are derived  
 29 from the generic functionalities of an observation data management system  
 30 given in the introduction.  
 31

- 32 1. The framework must support OLAP over large warehouses of spatial ob-  
 33 servation data.
- 34 2. The framework must support the definition of the behavior of *Internal*  
 35 *Processes* using a spatio-temporal declarative analysis language.
- 36 3. The language must support the integrated analysis of both conventional  
 37 E/R data and temporal, spatial and spatio-temporal sampled data.
- 38 4. The language must support aggregation functionality.  
 39

40 Based on the above requirements, first an XML Mapping Analysis Lan-  
 41 guage (MAPAL) is defined and next it is shown how it can be used for the  
 42 definition of internal analytical observation processes that are executed by  
 43 SODA during ETL.  
 44

### 45 4.1 Mapping Analysis Language

46  
 47 The use of data model based on functions and not on sets or sequences dis-  
 48 ables the direct use of some extension of well known languages like SQL and  
 49  
 50  
 51  
 52  
 53  
 54  
 55  
 56  
 57  
 58  
 59  
 60  
 61  
 62  
 63  
 64  
 65

XQuery. However, the hybrid logical-functional paradigm of the proposed MAPAL language reuses constructs of those well known languages and combines them with an XML syntax that simplifies their insertion in currently dominating web service interfaces. In particular, derived *Constants* and *Intensional* and *Extensional Mappings* are defined using three types of expressions, namely, *Functional Expressions*, *Conditional Expressions* and *Aggregate Expressions*. On the other hand, new *Dimensions* are defined using either *Sampling* or *Dimension Expressions*. The syntax and semantics of all those expressions are presented and illustrated with examples below.

*Functional Expressions*: A functional expression  $e$  defined in the context of variables  $v_1, v_2, \dots, v_n$ , denoted by  $e(v_1, v_2, \dots, v_n)$ , combines context variables with already defined constructors, system provided literals, operators, primitive mappings and type castings. The semantics are the obvious ones. The example below illustrates the use of functional expressions to define two *Constants* and a *Extensional Mapping*.

*Example 1* Obtain an extensional mapping that determines for each time instant whether vessel “Bur124” is navigating inside or outside its fishing temperature interval.

```
<Constant name="MinBur124">
  <Return>Vessel.MinFishingTemp("Bur124")</Return>
</Constant>

<Constant name="MaxBur124">
  <Return>Vessel.MaxFishingTemp("Bur124")</Return>
</Constant>

<ExtensionalMapping name="InsideTemperature" domain="GPS.Time t">
  <Return>
    SeaSurface.MODISTemperature(t,Vessel.Location(t,"Bur124"))&gt;=MinBur124
  AND SeaSurface.MODISTemperature(t,Vessel.Location(t,"Bur124"))&lt;=MaxBur124
  </Return>
</ExtensionalMapping>
```

Notice that functional expressions used to define constants do not have context variables. On the other hand, the functional expression used to define *Extensional Mapping* “InsideTemperature” has a context variable  $t$  whose scope is the dimension GPS.Time of its domain. Notice also that symbols  $<$  and  $>$  are not allowed in XML content and their entities  $&lt;$ ; and  $&gt;$ ; have to be used instead.

*Conditional Expressions*: They enable the introduction of if-then-else structures in the evaluation of *Constants*, *Intensional* and *Extensional Mappings*. The semantics are the obvious ones and the syntax is illustrated below with an example.

*Example 2* Define an Intensional Mapping SeaIce( $t, p$ ) that yields the risk of having ice in point  $p$  at time instant  $t$ . Risk “Red” for temperatures below -2 degrees, risk “Orange” for temperatures between 0 and -2 and risk “Green” for temperatures above 0 degrees.

```

1 <IntensionalMapping name="SeaIce" domain="t, p">
2 <When>fish:SeaSurface.MODISTemperature(t, p)&lt; -2 </When>
3 <ThenReturn>"Red"</ThenReturn>
4 <When>SeaSurface.MODISTemperature(t, p)&gt;= -2
5   AND SeaSurface.MODISTemperature(t, p)&lt;= 0
6 </When>
7 <ThenReturn>"Orange"</ThenReturn>
8 <ElseReturn>"Green"</ElseReturn>
9 </IntensionalMapping>

```

Set operations *Union* and *Intersection* between *Dimensions* are now formalized, as a prerequisite for the definition of the *Dimension* and *Aggregate Expressions* of MAPAL.

**Definition 9** Let  $d_1(T_1)$  and  $d_2(T_2)$  be two non-sampling *Dimensions*, where  $T_1$  and  $T_2$  are compatible data types, i.e., an implicit casting has been defined among them. Then  $d_1$  *Union*  $d_2$  is defined as the *Dimension*  $d(T) = \text{cast}(d_1 \text{ as } T) \cup \text{cast}(d_2 \text{ as } T)$ , where  $T$  is the result type of the implicit cast between  $T_1$  and  $T_2$  and  $\text{cast}(d_i \text{ as } T)$  is the *Dimension* resulting from casting each element of  $d_i$  to type  $T$ .

**Definition 10** Let  $d_1(T_1)$  and  $d_2(T_2)$  be two *Dimensions*, where  $T_1$  and  $T_2$  are compatible data types whose implicit result cast type is  $T$ , and at least one of the *Dimensions* is a *1D Sampling*. Then  $d_1$  *Union*  $d_2$  is defined as the *1D Sampling*  $S(m, M)$ , where

$$m = \min\{\text{cast}(vasT) | v \in d_1 \cup d_2\}$$

$$M = \max\{\text{cast}(vasT) | v \in d_1 \cup d_2\}$$

**Definition 11** Let  $d_1(T_1)$  and  $d_2(T_2)$  be two *Dimensions*, where  $T_1$  and  $T_2$  are compatible data types whose implicit result cast type is  $T$ , and at least one of the *Dimensions* is a *2D Sampling*. Then  $d_1$  *Union*  $d_2$  is defined as the *2D Sampling*  $S((x_m, y_m), (x_M, y_M))$ , where

$$x_m = \min\{\text{cast}(xasT) | (x, y) \in d_1 \cup d_2\}$$

$$y_m = \min\{\text{cast}(yasT) | (x, y) \in d_1 \cup d_2\}$$

$$x_M = \max\{\text{cast}(xasT) | (x, y) \in d_1 \cup d_2\}$$

$$y_M = \max\{\text{cast}(yasT) | (x, y) \in d_1 \cup d_2\}$$

**Definition 12** Let  $d_1(T_1)$  and  $d_2(T_2)$  be two *Dimensions* (either sampling or non-sampling), where  $T_1$  and  $T_2$  are compatible data types whose implicit cast type is  $T$ . Then  $d_1$  *Intersection*  $d_2$  is defined as the *Dimension*  $d(T) = \text{cast}(d_1 \text{ as } T) \cap \text{cast}(d_2 \text{ as } T)$ , where  $\text{cast}(d_i \text{ as } T)$  is the *Dimension* resulting from casting each element of  $d_i$  to type  $T$

Defining *Union* and *Intersection* in this way between samplings simplifies efficient implementation structures and matches the requirements of applications.

*Aggregate Expressions*: They are composed of the following four sections:

- *Variable scope specification*: It is composed of a sequence of variable scope specifications of the form

```

1  <ForEach var="v1">dimSetExpr1</ForEach>
2  <ForEach var="v2">dimSetExpr2</ForEach>
3  ...
4  <ForEach var="vn">dimSetExprN</ForEach>

```

where each  $vi$  is a variable name that iterates over the elements of the result of each *dimSetExpr<sub>i</sub>*. Each *dimSetExpr<sub>i</sub>* is a *Dimension Set Expression* of the form  $d1 \text{ OP } d2 \text{ OP } \dots \text{ OP } dm$ , where each  $di$  is a *Dimension* name and OP is either AND or OR. The semantics of AND and OR are those of *Dimension Set Operations Intersection* and *Union*, respectively.

- *Filtering*: The valid combinations of variable values may be restricted with an expression of the form

```
<Where>c(v1, v2, ..., vn)</Where>
```

where  $c(v1, v2, \dots, vn)$  is a *Functional Expression* that evaluates to boolean value true only for valid combinations of variable values. Valid combinations of variables are retrieved as a sequence of tuples, each of the form  $(v1, v2, \dots, vn)$ , ordered ascending by  $v1, v2, \dots, vn$ . A distinct ordering for the sequence may be optionally specified in the following *Ordering* section.

- *Ordering*: It is an optional sequence of expressions of either of the following two forms

```

<OrderAscendingBy>o(v1, v2, ..., vn)</OrderAscendingBy>
<OrderDescendingBy>o(v1, v2, ..., vn)</OrderDescendingBy>

```

where  $o$  is a functional expression whose result is of some ordered type. The tuple sequence resulting from the above sections is now ordered either ascending or descending by the result of these functional expressions.

- *Aggregate Evaluation*: It specifies an aggregate expression of the form

```
<Aggregate>a(v1, v2, ..., vn)</Aggregate>
```

where  $a$  combines functional expression elements with system provided aggregate functions. Variables  $vi$  must appear in  $a$  inside the scope of some aggregate function. Aggregate functions provided by the system may have the form of the classical ones of SQL (AVG, SUM, etc.) but they may also exploit the sequence ordering. Thus for example, function ATPOSITION( $S, n$ ) yields element located at position  $n$  in sequence  $S$ .

The following example illustrates the use of an aggregate expression for the definition of an *Extensional Mapping*.

*Example 3* Obtain the number of vessels that, at each time instant  $t$ , are navigating through areas that have some risk of ice.

*Dimension Expressions*: A new non-sampling *Dimension*  $d$  may be defined with an expression whose general form is as follows.

```

47 <Dimension name="d">
48   <ForEach var="v1">dimSetExpr1</ForEach>
49   <ForEach var="v2">dimSetExpr2</ForEach>
50   ...

```

51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

```

1  <ForEach var="vn">dimSetExprN</ForEach>
2  <Where>c(v1, v2, ..., vn)</Where>
3  <Return>e(v1, v2, ..., vn) </Return>
4 </Dimension>

```

The semantics are similar to those of *Aggregate Expressions*, however tuple sequence semantics are now replaced by tuple set semantics (as in the relational model).

*Example 4* Obtain a new dimension “boat3” that contains the names of the boats that at instant “2013-01-04T12:00:00” are navigating through a zone with temperature below 3 degrees.

```

13 <Dimension name="boat3">
14   <ForEach var="v">Vessel.Id</ForEach>
15   <Where>
16     SeaSurface.MODISTemperature("2013-01-04T12:00:00",
17       Vessel.Location("2013-01-04T12:00:00",v)) &lt; 3
18   </Where>
19   <Return>Vessel.Name(v)</Return>
20 </Dimension>

```

*Sampling Expressions:* A new *Sampling Dimension*  $S(m, M)$  may be defined with an expression of the form

```

24 <Dimension name="S">
25   <Start>m</Start>
26   <End>M</End>
27 </Dimension>

```

## 4.2 Definition of Analytical Processes

An XML language based on the MAPAL language described in the previous subsection enables the administrator of SODA to declaratively define internal analytical Observation Processes that are executed during ETL. Each internal process is of a specific Process Type, of those declared in XODDL. Therefore, during its definition, both a unique Process identifier and values for relevant Process Properties have to be provided. Besides, contrary to external processes, an internal process has a definition that is expressed with MAPAL. In particular, such a definition is composed of three sections.

1. *Preliminaries:* A possibly empty preliminary collection of temporary *Constant* and *IntensionalMapping* definitions that will be reused in the remainder sections.
2. *Time Dimension Definition:* Definition of the time *Dimension* of the result Process. It is recalled that each Process Type  $P$  of the framework has a *Dimension*  $P.Time$ , which is a *Sampling Dimension* in the case of time-triggered *Processes* and a non-sampling *Dimension* in the case of event-triggered *Processes*.

3. *Observed Property Definitions*: Definition of an *Extensional Mapping* with appropriate domain for each observed Feature Property whose values are generated by Processes of the present Process Type.

We shall focus here in points 2 and 3 above, since point 1 is just syntactic sugar. The time dimension of a Time-Triggered internal process type P of time resolution R is defined with an XML expression of the form

```
<TriggeredByTime>P1.Time, ..., Pn.Time</TriggeredByTime>
```

where each Pi.Time is the time *Dimension* of a Process Type Pi. The semantics are those of the *1D Sampling* S(m, M), where

$$m = \text{cast}(\min\{v | v \in P1.Time \cup \dots \cup Pn.Time\} \text{ as } \text{TimeInstant}(R))$$

$$M = \text{cast}(\max\{v | v \in P1.Time \cup \dots \cup Pn.Time\} \text{ as } \text{TimeInstant}(R))$$

The time dimension of an event-triggered internal process type P of time resolution R is defined with an XML expression of the form

```
<TriggeredByEvent>
  <Event var="t">P1.Time, P2.Time, ..., Pn.Time</Event>
  <Condition>c(t)</Condition>
</TriggeredByEvent>
```

where each Pi.Time is the time *Dimension* of a Process Type Pi and c(t) is a functional expression of boolean type. The semantics are those of the non-sampling *Dimension* defined by the set

$$\{\text{cast}(t \text{ as } \text{TimeInstant}(R)) | t \in P1.Time \cup \dots \cup Pn.Time \wedge c(t)\}.$$

Each feature property F.FP observed by an internal process type P has to be defined as an *Extensional Mapping* within the definition of P. The first dimension of the domain of F.FP will be P.Time. During ETL, each such *Extensional Mapping* will be evaluated, but restricting the scope of the evaluation to the P.Time dimension elements to be imported. This avoids the re-evaluation of the *Extensional Mapping* for the whole time extension of the data warehouse.

The definition of process types “IceAlert” and “FishingZone” of the running example are given next for illustration purposes.

```
<?xml version="1.0" encoding="utf-8"?>
<pd:ProcessDefinitions
  xmlns="es.usc.citius.de.mapal"
  xmlns:pd="es.usc.citius.de.soda.ProcessDefinition">
  <pd:Process id="IceAlert" processType="IceAlert">
    <pd:Definition>
      <IntensionalMapping name="IceRisk" domain="temperature">
        <When>temperature &lt;= -2</When><ThenReturn>"High"</ThenReturn>
        <When>temperature &lt; 0</When><ThenReturn>"Medium"</ThenReturn>
      </IntensionalMapping>
      <IntensionalMapping name="ExistsVesselInRisk" domain="t">
        <ForEach var="v">Vessel.Id</ForEach>
        <Where>Vessel.Temperature(t, v) &lt; 0</Where>
        <Aggregate>not EMPTY(v)</Aggregate>
      </IntensionalMapping>
    <pd:TriggeredByEvent>
      <pd:Event var="t">CTD.Time</pd:Event>
      <pd:Condition>ExistsVesselInRisk(t)</pd:Condition>
    </pd:TriggeredByEvent>
```

```

1      <ExtensionalMapping name="Vessel.IceAlert"
2          domain="IceAlert.Time t, Vessel.Id v">
3          <Return>IceRisk(Vessel.Temperature(t, v))</Return>
4      </ExtensionalMapping>
5  </pd:Definition>
6  </pd:Process>
7
8  <pd:Process id="FishingZone" processType="FishingZone">
9      <pd:Definition>
10         <pd:TriggeredByTime>MODIS.Time</pd:TriggeredByTime>
11         <ExtensionalMapping name="Vessel.FishingZone"
12             domain="FishingZone.Time t, Vessel.Id v">
13             <ForEach var="p">SeaSurface.Loc</ForEach>
14             <Where>SeaSurface.MODISTemperature(t, p) >=
15                 Vessel.MinFishingTemp(v)
16                 AND SeaSurface.MODISTemperature(t, p) <=
17                 Vessel.MaxFishingTemp(v)
18             </Where>
19             <Aggregate>VECTORIZE(p)</Aggregate>
20         </ExtensionalMapping>
21     </pd:Definition>
22     <Description>
23         This process obtains the piece of sea surface with
24         appropriate water temperature for fishing.
25     </Description>
26 </pd:Process>
27 </pd:ProcessDefinitions>

```

First “IceAlert” is defined as an event-triggered Process that is fired every time a CTD measures a temperature below 0. The Vessel.IceAlert Feature Property is then defined as either “High” or “Medium” depending on the temperature value. Notice that a vessel might measure a temperature above zero at the same time another one measures a temperature below zero. In that case, *Extensional Mapping* “Vessel.IceAlert” would return an undefined value for the first vessel. The definition of “FishingZone” Process illustrates the use of aggregate function VECTORIZE to generate a vector polygon from the set of Point2D elements of a 2D *Sampling* for which a specific condition holds.

## 5 Column-Oriented Implementation Issues

This sections discusses some issues that are guiding to the column-oriented implementation of SODA that is currently being undertaken.

### 5.1 SODA Architecture

A general overview of the architecture of the framework is depicted in Fig. 6. At the bottom of the figure, on the left side data catalog and system catalog are shown. The former records metadata related to *Feature* and *Process Types* and their relevant *Properties*. The latter records metadata of system provided constructors, including data types, primitive mappings and operators, aggregate functions and type castings. *Dimensions* and *Extensional Mappings* are recorded in one or various *Data Storage Units*. Besides, appropriate indexes are

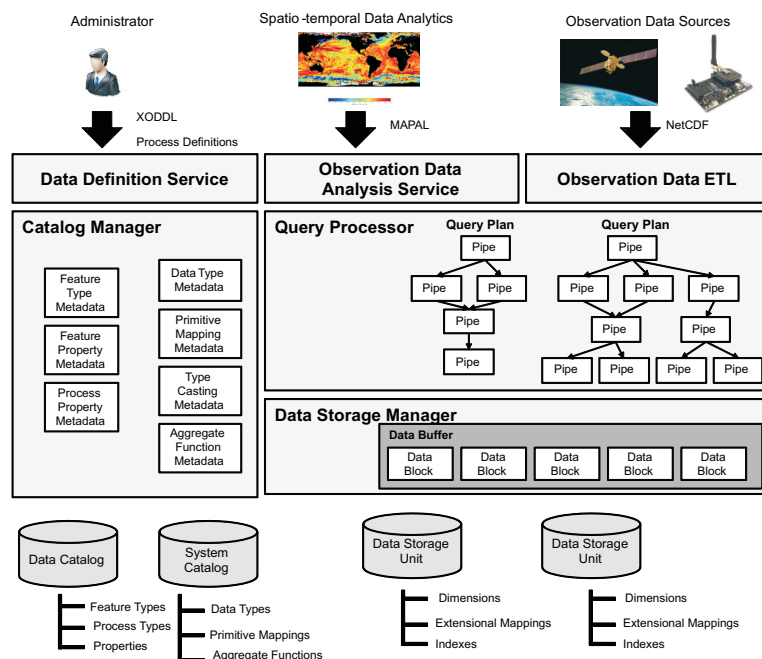


Fig. 6 Overview of SODA Architecture.

also recorded to speed up direct access to the above data structures. Following common practice in column-oriented DBMS implementations, lightweight compression techniques will be adopted for physical data representation (see Subsection 5.2). The *Data Storage Manager* enables efficient data access to the storage units, implementing buffering of data blocks. Access to all the system metadata is provided through the *Catalog Manager*. *Feature* and *Process* Metadata are generated from the schema definitions provided by the system administrator through the *Data Definition Service*. Besides the interpretation of XODDL (see Subsection 3.2) and internal process definitions (see Subsection 4.2), the *Data Definition Service* must enable the insertion of external processes together with their relevant *Process Property* values and non-observed feature data. Spatio-temporal data analytics are supported by an *Observation Data Analysis Service*, whose main functionality is the parsing and execution of MAPAL scripts. The execution of MAPAL query plans is performed by the *Query Processor*. Each query plan is represented by a *Direct Acyclic Graph* of pipes. Each pipe is generally executed in a distinct thread and supports a specific operation of the physical algebra described in Subsection 5.3. Finally, the *Observation Data ETL* component supports the execution of the internal analytical processes defined in the catalog during each observation data Extract Transform and Load (ETL).

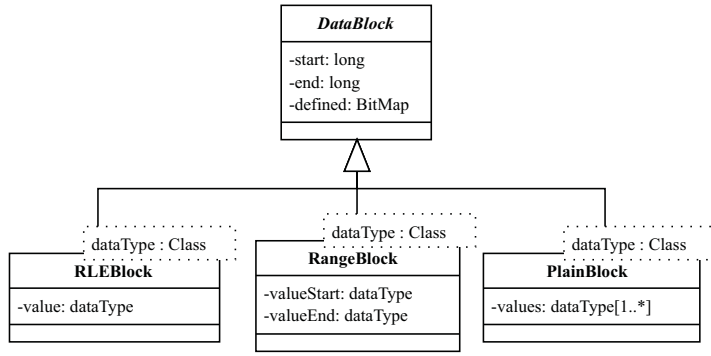


Fig. 7 Data Block Types in SODA.

## 5.2 Physical Data Representation

Various issues related to the physical data representation that are being considered during the current implementation are now discussed. Generally, each *Dimension* and *Extensional Mapping* is recorded as a header and a sequence of data blocks. *Dimension* data is recorded in self order whereas *Extensional Mappings* are recorded in the order of their respective *Dimensions*. Following common practices in column-oriented implementations, light weight compression techniques may be used to record both *Dimensions* and *Extensional Mappings*.

*Light Weight Compression:* The use of Light Weight Compression produces great performance improvement in current column-oriented implementations [5]. Thus, a requisite for the present SODA implementation is its flexibility to incorporate new compression techniques for both *Dimensions* and *Extensional Mappings*, enabling if possible the direct processing of compressed data. The current ongoing implementation uses three types of *DataBlock*, whose in-memory relevant data structures are depicted in Fig. 7 as a UML class diagram. A *RLEBlock* stores an element of a Run Length Encoding (RLE) sequence, i.e., it represents the repetition of a value from a start position to an end position. A *RangeBlock* represents a sequence of consecutive values of a data type by just recording the start and end values. A *PlainBlock* records a sequence of uncompressed values of a data type. the attribute “defined” of *DataBlock* records a compressed Bitmap that specifies which of the elements of *RLEBlocks* and *RangeBlocks* are valid.

A *Sampling*  $S(c_1, c_2)$  dimension is recorded as a single *RangeBlock*, that records the boundaries  $c_1$  and  $c_2$  of  $S$ . Non-sampling *Dimensions* are recorded with sequences of self ordered *PlainBlocks*, i.e., non-compressed data. In the future, delta encodings will be added to support efficient compression of numeric non-sampling *Dimensions*. Regarding *Extensional Mappings*, they are recorded as sequences of *DataBlocks*, ordered by the dimensions of their domain. Either *RLEBlocks* or *PlainBlocks* are used to record *Extensional Mappings*. Other

1 compression techniques including delta encodings and dictionary encodings  
2 will be added in the future. Properties *start* and *end* of a *DataBlock* are used  
3 in memory to reference the range of positions that represents the block inside  
4 its *Dimension* or *Extensional Mapping*. Property *defined* is used to determine  
5 which elements of the block are valid, i.e., which of them are distinct from *Un-*  
6 *defined*. These *defined* bitmaps may be set as a consequence of the evaluation  
7 of a where clause of MAPAL.  
8

9  
10 *Ordering*: As it was already stated, *Dimensions* are self ordered whereas *Ex-*  
11 *tensional Mappings* are ordered by the *Dimensions* of their domains. The time  
12 dimension of defined processes will always be the first *Dimension* of relevant  
13 *Extensional Mappings* of observed properties. This decision facilitates append-  
14 ing new temporal data without having to insert data values in the middle of  
15 already existing *Extensional Mappings*. *1D Samplings* are implicitly ordered  
16 from lower to higher values. Regarding *2D Samplings*, currently a linear order-  
17 ing is used, however, in the future other orderings defined by other Space  
18 Filling Curves shall also be considered.  
19

20  
21 *Standardized data interchange*: The interchange of data between the frame-  
22 work and its data providers and data consumer is based on the use of the  
23 NetCDF standard file format. This format enables the efficient representation  
24 of array data and it is broadly used in meteorological and oceanographic ap-  
25 plications. Currently it is also an OGC standard. Roughly speaking, data in  
26 a NetCDF is organized in dimensions and variables defined on those dimen-  
27 sions. These NetCDF concepts fit very well with the *Sampling* and *Extensional*  
28 *Mapping* concepts of the present framework.  
29

### 30 5.3 Query Processing

31  
32  
33 Issues related to the evaluation of MAPAL expressions are now discussed. In  
34 particular, first the physical algebra that is behind query processing is defined  
35 and next a couple of issues related to the current implementation of this algebra  
36 are discussed.  
37

#### 38 5.3.1 Physical Algebra

39  
40 The proposed physical algebra is a many-sorted algebra over elements of four  
41 different data structures, namely, *Dimensions*, *Domains*, *Constants* and *Ex-*  
42 *tensional Mappings*. *Dimensions*, *Constants* and *Extensional Mappings* have  
43 already been defined in Subsection 3.1. A *Domain* is defined as either a *Dimen-*  
44 *sion* or a Cartesian Product of 2 or more not necessarily distinct *Dimensions*.  
45 Operations of the algebra are classified into four different groups according  
46 to the result structure that they produce. In general, the signature of each  
47 operation has the form  
48

49  $OperatorName[ParameterList](ArgumentList)$   
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Table 1** Dimension and Domain Physical Algebra Operators

Dimension Operators	
$DimensionScan[dimName]()$	Parameter $dimName$ is the name of a <i>Dimension</i> . It generates a <i>RangeBlock</i> describing the range of positions in $dimName$ .
$Union(d_1, d_2)$	It performs the set <i>Union</i> of dimensions $d_1$ and $d_2$ (see Subsection 4.1).
$Intersection(d_1, d_2)$	It performs the set <i>Intersection</i> of dimensions $d_1$ and $d_2$ (see Subsection 4.1).
$Sampling(c_1, c_2)$	It produces a sampling dimension whose boundaries are defined by constants $c_1, c_2$ . The result is encoded as a <i>RangeBlock</i> .
$DupRem(M)$	It generates the new <i>Dimension</i> resulting from the elimination of duplicates from the codomain of <i>Extensional Mapping</i> $M$ .
Domain Operators	
$Domain(d)$	It generates a <i>Domain</i> of just one component from dimension $d$ .
$Product(D, d)$	It performs the Cartesian Product between <i>Domain</i> $D$ and <i>Dimension</i> $d$ .
$Select(M)$	If <i>Extensional Mapping</i> $M$ has a codomain of boolean type, then it selects from the <i>Domain</i> of $M$ the combinations where $M$ yields true.

where parameter list is a comma separated list of parameters and ArgumentList is a comma separated list of arguments. Arguments are denoted by different characters, possibly subscripted, depending on their type. Thus  $d_i, D_i, C_i$  and  $M_i$  denote respectively *Dimensions, Domains, Constants* and *Extensional Mappings*. Operations that produce *Dimensions* and *Domains* are described in Table 1. Table 2 describes operations that evaluate *Constants* and *Extensional Mappings*.

### 5.3.2 Pipeline Implementation

The physical algebra of the previous section is implemented following a *Producer-Driven Pipelining* approach. Thus, each operator is implemented as a pipe that is generally executed in a different thread. Data and reference interchange between pipes is asynchronous and it is supported by appropriate *DataBlock* buffers. This enables *Constant* and *Mapping* evaluation to be performed in parallel by different processing units, if they are available. Besides, if various storage units are available, different *Dimensions* and *Extensional Mappings* might be retrieved from storage also in parallel. This is known in the literature by pipeline parallelism [22]. The degree of pipeline parallelism is increased by the vertical partitioning approach followed by column-oriented implementations as the present one. An evaluation plan using the above operation pipes for the query of Example 3 is given in Fig. 8. First, two *DimensionScan* pipes

**Table 2** Constant and Mapping Physical Algebra Operators

<b>Constant Operators</b>	
<i>ConstantFetch</i> [consName]()	It obtains from storage the value recorded for <i>Constant consName</i>
<i>Literal</i> [literal]()	It produces the value represented by <i>literal</i> .
<i>IntensionalMapping</i> [iMap]( $C_1, C_2, \dots, C_n$ )	<i>iMap</i> is the name of either a primitive mapping, casting or operator provided by the system or an <i>Intensional Mapping</i> defined by the user. It evaluates <i>iMap</i> with arguments $C_1, C_2, \dots, C_n$ .
<i>ExtensionalMapping</i> [eMap]( $C_1, C_2, \dots, C_n$ )	<i>eMap</i> is the name of a <i>Extensional Mapping</i> . It obtains from storage the value of eMap for the domain element defined by $C_1, C_2, \dots, C_n$ .
<i>Conditional</i> ( $C_1, C_2, [C_3]$ )	If $C_1$ is true then it yields $C_2$ otherwise, if $C_3$ is specified then it yields $C_3$ , otherwise it yields <i>Undefined</i> .
<i>Aggregate</i> [aFun]( $M_1, M_2, \dots, M_n, [MO_1, MO_2, \dots, MO_m]$ )	All the argument <i>Extensional Mappings</i> are defined over the same <i>Domain</i> . It first produce a sequence of n-tuples from the codomains of $M_1, M_2, \dots, M_n$ , ordered by the codomains of $MO_1, MO_2, \dots, MO_m$ . Then it evaluates Aggregate Function AFun over such an ordered sequence.
<b>Mapping Operators</b>	
<i>Constant</i> ( $D, C$ )	It generates an <i>Extensional Mapping</i> that yields the result of evaluating <i>Constant C</i> for each element of Domain $D$ .
<i>Project</i> [dimRef]( $D$ )	For each element of Domain $D$ it yields the value of <i>Dimension</i> referenced by <i>dimRef</i> of $D$ .
<i>IntensionalMapping</i> [iMap]( $M_1, M_2, \dots, M_n$ )	<i>iMap</i> is the name of either a primitive mapping, casting or operator provided by the system or an <i>Intensional Mapping</i> defined by the user. It evaluates <i>iMap</i> using as arguments the values obtained from the evaluation of $M_1, M_2, \dots, M_n$ , for each element of their common <i>Domain</i> .
<i>ExtensionalMapping</i> [eMap]( $M_1, M_2, \dots, M_n$ )	<i>eMap</i> is the name of a <i>Extensional Mapping</i> . For each element of the common <i>Domain</i> of $M_1, M_2, \dots, M_n$ , it obtains from storage the value of <i>eMap</i> , using as domain the n-tuple resulting from the evaluation of $(M_1, M_2, \dots, M_n)$ .
<i>Conditional</i> ( $M_1, M_2, [M_3]$ )	For each element of the common <i>Domain</i> of $M_1, M_2$ and $M_3$ , if $M_1$ evaluates to true then it yields the result of evaluating $M_2$ otherwise, if $M_3$ is specified then it yields the result of evaluating $M_3$ , otherwise it yields <i>Undefined</i> .
<i>Aggregate</i> [groupBy][aFun]( $M_1, M_2, \dots, M_n, [MO_1, MO_2, \dots, MO_m]$ )	All the argument <i>Extensional Mappings</i> are defined over the same <i>Domain</i> . The sequence of n-tuples resulting from the evaluation of $M_1, M_2, \dots, M_n$ are grouped by dimension references <i>groupBy</i> . Each group of n-tuples is ordered by the result of evaluating $MO_1, MO_2, \dots, MO_m$ . Finally, Aggregate Function AFun is evaluated in the scope of each group.

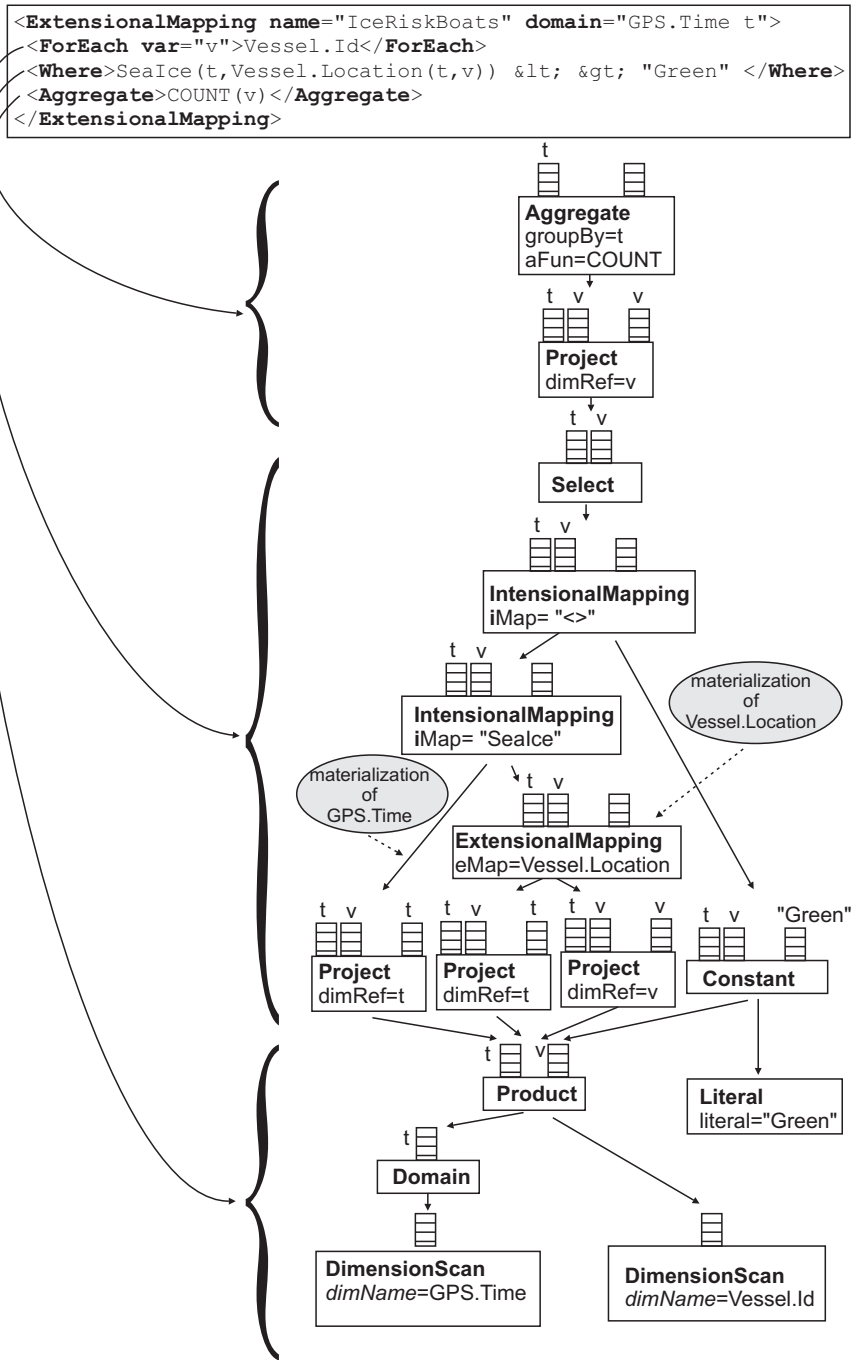


Fig. 8 Example of Query Evaluation Plan.

1 obtain the range of reference positions of the two involved dimensions. Notice  
2 that with just one *RangeBlock* for each pipe, all the reference positions of  
3 both *Dimensions* are represented. After passing through *Domain* and *Prod-*  
4 *uct* pipes we have the Cartesian Product of these two input *RangeBlocks*.  
5 Such a Cartesian Product may efficiently be represented with a combination  
6 of *RLEBlocks* for *Dimension t* and *RangeBlocks* for *Dimension v*. Next, three  
7 *Project* pipes generate three *Extensional Mappings* with identical *Domain* and  
8 again just reference positions in their codomains. Notice that up to this point,  
9 actual data values have not been read from storage yet. A *Constant* pipe gener-  
10 erates another *Extensional Mapping* that yields a constant value “Green” for  
11 each combination of GPS.Time and Vessel.Id. A *Extensional Mapping* pipe  
12 retrieves elements of Vessel.Location using as input the reference positions  
13 produced by the *Project* pipes. Given that GPS.Time and Vessel.Id form in  
14 that order the *Domain* of Vessel.Location, then position references do not have  
15 to be materialized to data in order to access Vessel.Location. Next, an *Inten-*  
16 *sional Mapping* pipe evaluates “SeaIce” for each combination of GPS.Time  
17 and Vessel.Location. Notice that now, GPS.Time reference positions have to  
18 be materialized to TimeInstant values in order to be able to evaluate “SeaIce”.  
19 Next, predicate “<>” is evaluated between the result of “SeaIce” and the *Con-*  
20 *stant* “Green”. At this stage, the result *Extensional Mapping* yields a boolean  
21 value for each combination of reference positions (t,v) from GPS.Time and  
22 Vessel.Id. Domain pipe *Select* uses such a boolean *Extensional Mapping* to re-  
23 strict the combinations to only those that are valid. Next, again a *Project* pipe  
24 is applied and finally an *Aggregate* pipe yields the expected result *Extensional*  
25 *Mapping* by grouping by *Dimension t* and counting valid reference positions  
26 of v.  
27

28 It is noticed that only the required data elements from *Dimensions* and  
29 *Extensional Mappings* are read from storage when they are required for some  
30 evaluation. Actually, in the example plan only values of GPS.Time *Dimen-*  
31 *sion* and Vessel.Location *Extensional Mapping* are obtained from disk. This  
32 approach of delaying data access operations as much as possible is called *Late*  
33 *Materialization* and it is a well known technique in column-oriented imple-  
34 mentations. It is also important to note that although the same *Extensional*  
35 *Mapping* or *Dimension* might appear various times in a evaluation plan, the  
36 appropriate use of buffers by the *Data Storage Manager* should avoid having  
37 to read their values from tertiary storage more than once. Specific structures  
38 might be required to achieve this as it has already been reported in [6].  
39  
40  
41  
42

## 43 6 Conclusions and Further Work

44 A framework for the analysis of spatial observation data was designed and  
45 qualitatively evaluated and compared with related data management tech-  
46 nologies and approaches. In particular, an observation data model was formal-  
47 ized based on the previous definition of a spatial data model. A declarative  
48 spatio-temporal data analysis language was also described and physical imple-  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 mentation issues were discussed. The advantages of the proposed framework  
2 may be summarized as follows.  
3

- 4 – Observation data semantics are directly supported by the spatial obser-  
5 vation data model and also incorporated in the declarative definition of  
6 *Internal Processes*.
- 7 – The framework provides direct support for the integrated representation  
8 and analysis of both conventional E/R data and temporal, spatial and  
9 spatio-temporal sampled data.
- 10 – The formalized temporal and spatial data types support the representation  
11 and transformation between different data resolutions.
- 12 – Both data (*Extensional Mappings*) and behavior (*Intensional Mappings*)  
13 are represented with the well known mathematical concept of function.  
14 Therefore, it is likely that the approach will be friendly to scientific users.  
15 Besides, a functional approach simplifies the definition and reuse of inter-  
16 mediate results, enabling the well known software engineering Black Box  
17 concept.
- 18 – The use of XML syntax of the proposed languages (XODDL and MAPAL)  
19 simplifies their incorporation in web service interfaces and enables the use  
20 of widely adopted XML technologies in the implementation.
- 21 – It is estimated that the use of a single non-nested data structure in the  
22 data model will simplify the efficient implementation of the framework.  
23

24  
25 The main drawback of SODA arise from the assumption of a brand new  
26 data management paradigm departing from the classical relational-SQL one,  
27 which is a clear handicap for current DBMS users. However, this is alleviated  
28 by the fact that well known logical and functional formalisms were combined to  
29 define MAPAL, which makes its constructors very similar to those of currently  
30 available languages like XQuery.

31 Future work issues include the following. Complete a first prototype of the  
32 framework in a single node multi-core hardware architecture and test its per-  
33 formance in comparison with other relevant solutions; Incorporate query op-  
34 timization techniques and appropriate indexing structures to improve system  
35 performance; Redesign and implement a new version for a multi-node share  
36 nothing architecture, exploiting the parallelism with appropriate horizontal  
37 partitioning of *Dimensions* and *Extensional Mappings*.  
38

39  
40 **Acknowledgements** This work has been partially supported by the Spanish Ministry  
41 of Science and Innovation (TIN2010-21246-C02-02). The authors are also grateful to the  
42 reviewers, whose comments contributed to greatly improve the paper.  
43

## 44 References

- 45 1. Apache cassandra (2014). URL <http://cassandra.apache.org/>
  - 46 2. MongoDB (2014). URL <http://www.mongodb.org/>
  - 47 3. Vertica (2014). URL <http://www.vertica.com/>
  - 48 4. VoltDB (2014). URL <http://voldb.com/>
- 49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1 5. Abadi, D., Madden, S., Ferreira, M.: Integrating compression and execution in  
2 column-oriented database systems. In: Proceedings of the 2006 ACM SIG-  
3 MOD International Conference on Management of Data, SIGMOD '06, pp. 671–  
4 682. ACM, New York, NY, USA (2006). DOI 10.1145/1142473.1142548. URL  
5 <http://doi.acm.org/10.1145/1142473.1142548>
- 6 6. Abadi, D., Myers, D., DeWitt, D., Madden, S.: Materialization strategies in a column-  
7 oriented dbms. In: Data Engineering, 2007. ICDE 2007. IEEE 23rd International Con-  
8 ference on, pp. 466–475 (2007). DOI 10.1109/ICDE.2007.367892
- 9 7. Arasu, A., Babu, S., Widom, J.: The cql continuous query language: semantic founda-  
10 tions and query execution. *The VLDB Journal* **15**(2), 121–142 (2006). DOI  
11 10.1007/s00778-004-0147-z. URL <http://dx.doi.org/10.1007/s00778-004-0147-z>
- 12 8. Baumann, P., Dehmel, A., Furtado, P., Ritsch, R., Widmann, N.: The multidim-  
13 ensional database system radsaman. In: Proceedings of the 1998 ACM SIG-  
14 MOD international conference on Management of data, SIGMOD '98, pp. 575–  
15 577. ACM, New York, NY, USA (1998). DOI 10.1145/276304.276386. URL  
16 <http://doi.acm.org/10.1145/276304.276386>
- 17 9. Baumann, P., Holsten, S.: A comparative analysis of array models for databases.  
18 In: T.h. Kim, H. Adeli, A. Cuzzocrea, T. Arslan, Y. Zhang, J. Ma, K.i. Chung,  
19 S. Mariyam, X. Song (eds.) *Database Theory and Application, Bio-Science and Bio-*  
20 *Technology, Communications in Computer and Information Science*, vol. 258, pp. 80–  
21 89. Springer Berlin Heidelberg (2011). DOI 10.1007/978-3-642-27157-19. URL  
22 <http://dx.doi.org/10.1007/978-3-642-27157-19>
- 23 10. Bowers, S., Madin, J., Schildhauer, M.: A conceptual modeling framework for expressing  
24 observational data semantics. In: Q. Li, S. Spaccapietra, E. Yu, A. Oliv (eds.) *Concept-*  
25 *tual Modeling - ER 2008, Lecture Notes in Computer Science*, vol. 5231, pp. 41–54.  
26 Springer Berlin Heidelberg (2008). DOI 10.1007/978-3-540-87877-35. URL  
27 <http://dx.doi.org/10.1007/978-3-540-87877-35>
- 28 11. Bröring, A., Stasch, C., Echterhoff, J.: OGC Sensor Observation Ser-  
29 vice Interface Standard. Open Geospatial Consortium (OGC) (2012).  
30 [Http://www.opengeospatial.org/standards/sos](http://www.opengeospatial.org/standards/sos)
- 31 12. Brown, P.G.: Overview of scidb: large scale array storage, processing and analysis.  
32 In: Proceedings of the 2010 ACM SIGMOD International Conference on Manage-  
33 ment of data, SIGMOD '10, pp. 963–968. ACM, New York, NY, USA (2010). DOI  
34 10.1145/1807167.1807271. URL <http://doi.acm.org/10.1145/1807167.1807271>
- 35 13. Cerveira Cordeiro, J.a.P., Câmara, G., Moura De Freitas, U., Almeida, F.: Yet another  
36 map algebra. *Geoinformatica* **13**(2), 183–202 (2009). DOI 10.1007/s10707-008-0045-4.  
37 URL <http://dx.doi.org/10.1007/s10707-008-0045-4>
- 38 14. Compton, M., Barnaghi, P., Bermudez, L., Garca-Castro, R., Corcho, O., Cox, S., Gray-  
39 beal, J., Hauswirth, M., Henson, C., Herzog, A., Huang, V., Janowicz, K., Kelsey, W.D.,  
40 Phuoc, D.L., Lefort, L., Leggieri, M., Neuhaus, H., Nikolov, A., Page, K., Passant, A.,  
41 Sheth, A., Taylor, K.: The {SSN} ontology of the {W3C} semantic sensor network in-  
42 cubator group. *Web Semantics: Science, Services and Agents on the World Wide Web*  
43 **17**(0), 25 – 32 (2012). DOI <http://dx.doi.org/10.1016/j.websem.2012.05.003>. URL  
44 <http://www.sciencedirect.com/science/article/pii/S1570826812000571>
- 45 15. Cox, S.: Geographic Information - Observations and Measurements. Open Geospa-  
46 tial Consortium (OGC) Abstract Specification Topic 20 and ISO 19156:2011(E) (2013).  
47 [Http://www.opengeospatial.org/standards/om](http://www.opengeospatial.org/standards/om)
- 48 16. Cugola, G., Margara, A.: Processing flows of information: From data stream to  
49 complex event processing. *ACM Comput. Surv.* **44**(3), 15:1–15:62 (2012). DOI  
50 10.1145/2187671.2187677. URL <http://doi.acm.org/10.1145/2187671.2187677>
- 51 17. Date, C.J., Darwen, H., Darwen, H.: *Temporal Data and the Relational Model: A De-*  
52 *tailed Investigation into the Application of Interval and Relation Theory to the Problem*  
53 *of Temporal ... Kaufmann Series in Data Management Systems*, 1 edn. Morgan Kauf-  
54 man Publ Inc (2002)
- 55 18. Galpin, I., Brenninkmeijer, C., Gray, A., Jabeen, F., Fernandes, A., Paton, N.:  
56 Snee: a query processor for wireless sensor networks. *Distributed and Paral-*  
57 *lel Databases* **29**(1-2), 31–85 (2011). DOI 10.1007/s10619-010-7074-3. URL  
58 <http://dx.doi.org/10.1007/s10619-010-7074-3>

- 1 19. Gray, P.M.D.: *The Functional Approach to Data Management: Modeling, Analyzing,*  
2 *and Integrating Heterogeneous Data.* SpringerVerlag (2004)
- 3 20. Güting, R.H.: *Spatial Databases.* John Wiley & Sons, Inc. (2001). DOI  
4 10.1002/047134608X.W4317. URL <http://dx.doi.org/10.1002/047134608X.W4317>
- 5 21. Güting, R.H., Böhlen, M.H., Erwig, M., Jensen, C.S., Lorentzos, N.A., Schneider,  
6 M., Vazirgiannis, M.: A foundation for representing and querying moving objects.  
7 *ACM Trans. Database Syst.* **25**(1), 1–42 (2000). DOI 10.1145/352958.352963. URL  
8 <http://doi.acm.org/10.1145/352958.352963>
- 9 22. Harizopoulos, S., Shkapenyuk, V., Ailamaki, A.: Qpipe: A simultaneously pipelined  
10 relational query engine. In: *Proceedings of the 2005 ACM SIGMOD Inter-*  
11 *national Conference on Management of Data, SIGMOD '05*, pp. 383–394.  
12 ACM, New York, NY, USA (2005). DOI 10.1145/1066157.1066201. URL  
13 <http://doi.acm.org/10.1145/1066157.1066201>
- 14 23. Idreos, S., Groffen, F.E., Nes, N.J., Manegold, S., Mullender, K.S., Kersten,  
15 M.L.: MonetDB: Two Decades Of Research In Column-Oriented Database Ar-  
16 chitectures. *IEEE Data Engineering Bulletin* **35**(1), 40 – 45 (2012). URL  
17 <http://oai.cwi.nl/oai/asset/19929/19929B.pdf>
- 18 24. International Organization for Standardization (ISO): *Information technology –*  
19 *Database languages – SQL multimedia and application packages – Part 3: Spatial.*  
20 *ISO/IEC 13249-3:2011* (2011)
- 21 25. Jain, N., Mishra, S., Srinivasan, A., Gehrke, J., Widom, J., Balakrishnan, H.,  
22 Çetintemel, U., Cherniack, M., Tibbetts, R., Zdonik, S.: Towards a stream-  
23 ing sql standard. *Proc. VLDB Endow.* **1**(2), 1379–1390 (2008). URL  
24 <http://dl.acm.org/citation.cfm?id=1454159.1454179>
- 25 26. Kulkarni, K., Michels, J.E.: Temporal features in sql:2011. *SIG-*  
26 *MOD Rec.* **41**(3), 34–43 (2012). DOI 10.1145/2380776.2380786. URL  
27 <http://doi.acm.org/10.1145/2380776.2380786>
- 28 27. Lorentzos, N.A., Viqueira, J.R.R.: Relational formalism for the management of spatial  
29 data. *The Computer Journal* **49**(1), 62–81 (2006). DOI 10.1093/comjnl/bxh136. URL  
30 <http://comjnl.oxfordjournals.org/content/49/1/62.abstract>
- 31 28. Madden, S.R., Franklin, M.J., Hellerstein, J.M., Hong, W.: Tinydb: an ac-  
32 quisitional query processing system for sensor networks. *ACM Trans.*  
33 *Database Syst.* **30**(1), 122–173 (2005). DOI 10.1145/1061318.1061322. URL  
34 <http://doi.acm.org/10.1145/1061318.1061322>
- 35 29. Madin, J., Bowers, S., Schildhauer, M., Krivov, S., Pennington, D., Villa, F.: An on-  
36 tology for describing and synthesizing ecological observation data. *Ecological Informatics* **2**(3), 279 – 296 (2007). DOI <http://dx.doi.org/10.1016/j.ecoinf.2007.05.004>.  
37 URL <http://www.sciencedirect.com/science/article/pii/S1574954107000362>. Meta-  
38 information systems and ontologies. A Special Feature from the 5th International Con-  
39 ference on Ecological Informatics ISEI5, Santa Barbara, CA, Dec. 47, 2006 Novel Con-  
40 cepts of Ecological Data Management S.I.
- 41 30. Neteler, M., Mitasova, H.: *Open Source GIS: A GRASS GIS Approach.* Third edition.  
42 Springer, New York, USA (2008)
- 43 31. Obe, R., Hsu, L.: *PostGIS in Action.* Manning Publications Co, Stamford, CT, USA  
44 (2011)
- 45 32. Open Geospatial Consortium (OGC): *OpenGIS Sensor Model*  
46 *Language (SensorML) Implementation Specification* (2007).  
47 [Http://www.opengeospatial.org/standards/sensorml](http://www.opengeospatial.org/standards/sensorml)
- 48 33. Sagan, H.: *Space-Filling Curves.* Springer (1994)
- 49 34. Schut, P.: *OpenGIS Web Processing Service.* Open Geospatial Consortium (OGC)  
50 (2007). [Http://www.opengeospatial.org/standards/wps](http://www.opengeospatial.org/standards/wps)
- 51 35. Snodgrass, R.T. (ed.): *The TSQL2 Temporal Query Language.* Kluwer (1995)
- 52 36. Stonebraker, M., Abadi, D.J., Batkin, A., Chen, X., Cherniack, M., Ferreira, M., Lau,  
53 E., Lin, A., Madden, S., O’Neil, E., O’Neil, P., Rasin, A., Tran, N., Zdonik, S.: C-  
54 store: A column-oriented dbms. In: *Proceedings of the 31st International Conference*  
55 *on Very Large Data Bases, VLDB '05*, pp. 553–564. VLDB Endowment (2005). URL  
56 <http://dl.acm.org/citation.cfm?id=1083592.1083658>

- 1 37. Vaisman, A., Zimányi, E.: A multidimensional model representing continuous fields  
2 in spatial data warehouses. In: Proceedings of the 17th ACM SIGSPATIAL Inter-  
3 national Conference on Advances in Geographic Information Systems, GIS '09, pp.  
4 168–177. ACM, New York, NY, USA (2009). DOI 10.1145/1653771.1653797. URL  
5 <http://doi.acm.org/10.1145/1653771.1653797>
- 6 38. Viqueira, J., Lorentzos, N.: Sql extension for spatio-temporal data. *The VLDB Journal*  
7 **16**(2), 179–200 (2007)
- 8 39. Zhang, Y., Kersten, M.L., Manegold, S.: SciQL: Array Data Processing In-  
9 side An RDBMS. In: Proceedings of ACM SIGMOD International Confer-  
10 ence on Management of Data 2013, pp. 1049 – 1052. ACM (2013). URL  
11 <http://oai.cwi.nl/oai/asset/21401/21401A.pdf>

11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65