

# Using Versions in GIS

Claudia Bauzer Medeiros<sup>1</sup>, Geneviève Jomier<sup>2</sup>

<sup>1</sup> DCC – UNICAMP, 13081-970 Campinas – SP – Brazil  
cmbm@dcc.unicamp.br

<sup>2</sup> LAMSADE – Université Paris IX Dauphine, 75015 Paris - France  
jomier@lamsade.dauphine.fr

**Abstract.** GIS have become important tools in public planning activities (e.g, in environmental or utility management). This type of activity requires the creation and management of alternative scenarios, as well as analysis of temporal data evolution. Existing systems provide limited support for these operations, and appropriate tools are yet to be developed. This paper presents a solution to this problem. This solution is based on managing temporal data and alternatives using the DBV version mechanism.

**Keywords:** versions, GIS, georeferenced data evolution.

## 1 Introduction

Geographic Information Systems (GIS) are automated systems that manipulate *georeferenced data* – data about geographic phenomena associated with their spatial relationships and location on the terrestrial surface. The term *georeferenced entity*, in this paper, refers to any type of entity whose components are georeferenced (be it a house, a residential plot or a geographic region).

GIS data contains *spatial* and *thematic* components. Thematic components are alphanumeric values related to georeferenced entities, e.g., the name of a mountain, or the type of vegetation cover. Texts or images are considered to be unconventional thematic data. Spatial data has two different properties: *geometric properties* such as size and shape of spatial objects; and *topological properties* such as adjacency and containment. Data are stored in *thematic layers*, which are combined in different ways in order to process a query.

GIS queries involve one or more of the following issues [Flo91, Peu93]:

1. *What* kind of phenomenon is this? (describe non-spatial characteristics)
2. *Where* is this phenomenon located? (describe the spatial characteristics of a given georeferenced entity, which comprises its location, topological and geometrical characteristics, and often its relationships with other entities).
3. *When* was this data collected? (determine data validity period).
4. What did this entity look like at some *past* period? What will happen to it in some *future* period? (examine previous states of an entity and predict its future evolution, given its recorded behavior).
5. *What* would happen to an entity *if* certain events were to take place? (simulation and comparison of alternative scenarios by changing existing data).

Database research for GIS support has been centered on *spatial databases*, especially in what concerns algorithms for storing and accessing spatial data (e.g., [Sam89]). Database researchers are concerned with providing users with fast access to georeferenced data by means of spatial indices. They are, furthermore, developing different types of query languages and mechanisms in order to allow processing the first two kinds of queries (i.e., *what*, *where*). Existing systems also provide different facilities for handling these two questions.

The remaining kinds of queries, however, involve other types of knowledge and distinct storage and indexing facilities. *When* processing (queries 3 and 4) requires temporal database management, itself a matter of intensive research (see, for instance, the discussion of open issues in [Sno92]).

The simulation of scenarios (query 5) is supported by some systems for specific situations, in a limited scale, using controlled parametrization of data values (e.g., [HGW93]). The combination of simulation results is, however, not allowed, especially when the user wants to compare alternatives. Users have to store the different scenarios in separate files, and have to handle themselves the management of these files, by embedding appropriate code into their applications.

This paper extends the work of [MJ93], presenting and detailing a framework which allows processing queries 1 through 5. This solution is based on the notion of database versions, and adopts the DBV version model of [CVJ94], now being simulated on the O2 database system.

The DBV model enables the simultaneous management of distinct alternative scenarios, as well as handling georeferenced feature evolution through time with considerable savings in space. The DBV model can be embedded in any database system and does not require that the end-user control the different data configurations.

This paper's main contributions are:

- analysis of the problem of georeferenced data evolution from a database point of view, rephrasing this problem in a versioning framework;
- description of how the DBV version model can support the management of this evolution, using concrete examples. So far, GIS have not considered version mechanisms, given the complexity of the factors involved.

The research discussed in this paper is part of the DOMUS[PMB93] environmental planning project, which uses real-life data from non-settled areas in the state of São Paulo, Brazil.

The rest of this paper is organized as follows. Section 2 characterizes GIS applications. Section 3 gives an overview of the DBV mechanism. Section 4 presents a detailed example of how this mechanism can support management of scenarios and data evolution. Section 5 shows how other version mechanisms are unable to perform this task satisfactorily. Section 6 presents conclusions and directions for future research.

## 2 GIS application requirements

GIS demand that DBMS keep track of massive amounts of georeferenced data, of different natures, collected using heterogeneous devices, and at different time periods.

The rapid growth in GIS has resulted in a large number of systems, each of which with its own data storage and handling characteristics. Present systems are based on relational database management systems (DBMS). However, several of the requirements of GIS applications are not provided by standard relational DBMS. Thus, special data handlers have been developed to interact with the stored relations and allow the management and display of georeferenced data.

Relational databases do not provide an adequate underlying model to support most types of geographic data [Ege92]: the use of tables with fixed number of attributes does not allow flexibility in the management of georeferenced data, nor the incremental development of applications. Thus, researchers have directed their attention to new architectures: extensible relational, object-oriented or rule-based systems [MP94].

In all of these, there is no consideration for temporal queries or comparison of scenarios (such as queries 3 through 5 of section 1). The difficulties posed to answering these queries involve factors that cannot be handled adequately by present GIS. The first factor is due to the nature of GIS data, which requires special indexing and buffering techniques. This is aggravated by the introduction of the time element, not supported in commercial databases. Thus, users are forced to manage time values themselves, if they want to analyze data evolution.

Another issue concerns data integration over time. Georeferenced data may be collected at different time periods. This creates another type of value inconsistency, which is due to the temporal evolution of georeferenced entities. Thus, if the mapping of a region takes several months (or sometimes years), differences will occur which must be taken into consideration. As pointed out by [Flo91], another related problem is the difference in time scales. Some phenomena (e.g., vegetation) fluctuate according to a seasonal cycle, whereas others (e.g, temperature) may vary on a daily basis. These factors must be considered for queries on the evolution of georeferenced phenomena for a given region.

## 3 The DBV Version Mechanism

As seen above, present GIS still lack facilities for providing the following services:

- automatic representation and management of data evolution in time;
- management of alternative scenarios for planning purposes.

These are the same type of problems that are faced by version mechanisms (even though the latter have not yet considered georeferenced data). Thus, it is only natural to examine the feasibility of adopting versions to allow such services.

### 3.1 Versions in databases

Versions are a means of storing different states of a given entity, thereby allowing the control of alternatives and of temporal data evolution. The management of versions in databases has centered on different ways for keeping files. Research has appeared mostly in the context of CASE systems and CAD/CAM projects (e.g., [KSW86, Kat90, BBA91, TG92, KS92, LST93]). The subjects discussed cover the creation and manipulation of entity versions, their identification, the handling of time, status, authorization, and concurrency mechanisms. Versions are also commonly proposed for dealing with concurrency control, especially for long transactions. In this last context, different users are granted access to copies (i.e., versions) of the same set of data, thus allowing them to work in parallel.

### 3.2 An overview of the DBV mechanism

Version mechanisms must put together entity versions to reflect a given state of the modelled universe. Existing approaches support this by means of chains of pointers, which keep track of connections among versions of a given entity, as well as among entities that belong to a given version state. Thus, the database is perceived as a set of entities connected by several linked chains. There is often confusion between version management and the underlying data model.

The DBV mechanism [CJ90, CVJ94] has a different approach. The main principles of this model are the following:

- The user should always be able to manage entity versions within a specific, identifiable, context. This context is called *database version* – *Dv* for short. *Dv* represents a “version”, or identifiable state, of the modelled universe.
- Each *Dv* contains a logical version of each (identified) entity of the database. The value of this version may be *nil*, meaning that it does not exist in a given *Dv*.
- If several logical versions of an entity have the same value, there is no physical replication. Rather, these logical versions share the same physical version. The mechanism of mapping logical/physical versions of entities is transparent to the user.

Thus, rather than keeping track of versions of individual entities, the problem is treated from a point of view where the unit of versioning is the *Dv* context, which corresponds to a state of the universe modelled by the database, regardless of the underlying data model. The notion of version as seen by the user (the *logical database versions*) is independent of what is actually stored (*physical versioning*), and there is no replication of data.

From a logical point of view, the database is perceived as a set of consistent database versions *Dv*, which can evolve independently of each other. The user works in this augmented database by selecting the *Dv* context(s) of interest. Once the desired *Dv* is selected, the user can treat it as a database on its own, querying and updating it. Logical operations on *Dv* are translated into actual

physical versioning operations by the version management system. Temporal and alternative data are naturally managed by this model.

The versioning of an entity is logically performed by versioning the entire logical database to which it belongs. The logical independence of  $Dv$  allows defining two types of update transactions: those that manipulate some  $Dv$ ; and those that derive a new  $Dv$ . The latter can be described in two steps. Let  $\mathcal{E}$  be an entity of a given  $Dv$ , for which the user wants to create a new version  $\mathcal{E}_1$ :

- first, a (logical) copy of  $Dv$  is created, corresponding to a new  $Dv_1$  database version, identical to  $Dv$ ;
- second,  $Dv_1$  is updated (i.e.,  $\mathcal{E}_1$  is created from  $\mathcal{E}$ , and additional updates are performed in order to maintain the consistency inside  $Dv_1$ ).

Thus, a multiversion database is a set of logical *consistent* states, each of which corresponds to a different version of the world, created by the user. Each  $Dv$  is, therefore, a unit of consistency and can evolve independently.

Physically, in order to properly associate an entity with its versions, the DBV mechanism relies on the notion of *identity*: each  $Dv$  has an associated identifier which is used for managing purposes. The identifier mechanism associates every logical entity version  $\mathcal{E}_i$  to its proper context  $Dv_i$ .

The retrieval of each  $Dv$  is automatically ensured by the version manager by examining tables of identifiers. Implementation details appear in [CJ90].

In this context, time may be supported by mapping timestamps into the DBV identifiers. Thus, temporal queries do not require handling of special attributes; rather, they are processed by the versioning mechanism, which puts together data that belongs to the same identifiable temporal state.

## 4 Applying the DBV model to GIS

This section shows an extended example of how the DBV approach allows handling GIS queries for evolution of data and comparison of alternatives.

Consider the following sequence of queries:

- Analysis of temporal data evolution: “What has been the observed modifications of the forest cover in a given area for a specific time period  $\mathcal{P}$ ?”
- Prevision of future based on recorded past: “What is the probability of this forest cover decreasing in area, given information collected along  $\mathcal{P}$ ? What are the possible damages – extent and intensity?”
- Comparative analysis between actual data and simulated scenarios: “Given the actual state of the forest cover in the area, how accurate were simulations performed along  $\mathcal{P}$  to determine its evolution during the same period?”

All these queries concern the same type of theme – vegetation – for the same geographical area. They require searching through different (historical) scenarios, first identifying the area and then its forest cover.

The first query requires doing a statistical analysis of a historical sequence of vegetation data. The second query requires using the previous analysis to

perform prediction of phenomena. The third query assumes that the two other queries were periodically posed in the past, and that their results were stored. Thus, it demands comparison of predicted behavior and actual observed behavior for several periods in the database's history.

In the DBV approach, the user is provided with an integrated view of the world: it is perceived as a set of database states ( $Dv$ ), each corresponding to an independent consistent version of the user's universe. The evolution of georeferenced entities is accompanied by the corresponding evolution of  $Dv$  states. For the user, there is no predefined link between different  $Dv$ , which allows working either on one context or navigating across contexts. This corresponds to what the user manipulates in reality, since no georeferenced phenomenon can be treated in isolation.

There is no difference in treatment for actual (measured) data values or (alternative) predicted values. Thus, several  $Dv$  may exist for the same time period, each describing a state of the world – either an actual recorded state or some alternative artificial state generated for planning purposes. Therefore, for any time period, the database may contain a set of logical databases: the modelled real world and different simulated scenarios. For vegetation cover, for instance, one can keep track of several parallel scenarios by modifying distinct parameters, e.g., rainfall or evolution of human settlements.

Thus, for any of the three queries, the processing is performed as:

- (i) Select all  $Dv$  within the specified time period;  
(The version manager performs this operation by accessing the identifier of each  $Dv$ , which in this case will contain a timestamp identification)
- (ii) For each such  $Dv$ , select the area and its cover, by performing standard GIS (nontemporal) database queries;  
(Each  $Dv$  selected in the previous step is seen by the user as an independent consistent database. Thus, it can be queried independently of the rest of the DBV database, regardless of other existing versions.)
- (iii) Perform the simulation operations on the set of areas and covers obtained from the execution of the two previous steps.

Physically, the creation of a logical database does not require physical duplication of entities, just creation of identifiers and recording of data changes (differential information). These changes and the corresponding identifiers are used to build the complete (logical) database states  $Dv$ .

The  $Dv$  are built by gathering together all entities present at a given database state, by means of special index structures. This means retrieving all entities that have compatible identifier values.

Query processing may be speeded up by using the notion of configuration. In fact, configurations characterize units of work inside a context (and thus of consistency). Thus, users may decide to specify a configuration containing the area and its cover (requested in the three queries). This will speed up version processing for this type of query (in a way similar to precomputed views).

## 5 Other version mechanisms and GIS

A good introduction to the problems of handling spatio-temporal data in GIS are the set of papers in [FCF92], which cover different issues. They range from problems in database support of time [Sno92] to discussing the concept of a region in creating study scenarios [Gut92]. The need for flexible mechanisms to allow managing of these scenarios is stressed in several papers.

Database research on versions has not dealt with problems related to GIS, and there are very few reports of GIS using version mechanisms (e.g., [Bat92, NTE92]). The main reasons for this are:

- Most version management mechanisms available in database systems become cumbersome when it comes to managing the evolution of instances. They require maintaining complex data structures to follow data evolution.
- GIS data is complex and occupies considerable space. Thus, its management already presents so many challenges to a DBMS that there is no question of coupling it to the usual version mechanisms.
- In many cases, it is impossible to follow the evolution of phenomena across time periods, since entities may disappear or suffer unexpected modifications.

Thus, even though versioning solves users' problems, it has not, so far, been seriously considered in the GIS context, and is used at most to support parallel access to data. Their use as a means to manage temporal data evolution is not considered. Rather, researchers consider (historical) versions of entire files – e.g., the sequence of file versions for a given thematic layer. Finally, when versions are associated with georeferenced entities, there is no concern with how to manage them from a database point of view.

The GFIS [Bat92] system uses a standard relational DBMS coupled to a geographic data manager. Version management is left to the database system, and is geared towards controlling parallel access. There is no possibility of selecting versions for queries, or of handling sequences of past states.

[NTE92] discuss different data structures for implementing versions on top of tables using an object-oriented language. The paper provides a comparative analysis of these structures, but does not apply them to real data.

We now briefly review how traditional version schemes would cope with the queries of section 4. In such schemes, two types of approach are possible:

• **Snapshot view** The complete data files (layers, with spatial and thematic information) have to be stored, together with time stamp indication. (I.e., the database is in fact a set of database snapshots, where each snapshot contains several thematic layers.) Thus, in order to answer the first two queries, the system has to:

- (i) retrieve the entire layer files for the time period;
- (ii) for each layer, select the desired area and determine its forest cover;
- (iii) produce the time series analysis desired.

In order to answer the third query, the database must contain not only the entire layers for every period of interest, but also files with simulation results.

The snapshot view, though relatively simple to process, entails massive storage occupation, and is therefore not feasible for practical purposes. Furthermore, it requires the actual recording of the entire database, and thus the variety and periodicity of stored phenomena is limited, due to size constraints.

• **Historical chain view** The history of entities is maintained through a linked list of data values and timestamps: only differential values are kept. (The database is seen as a conglomerate of linked chains in all directions.)

In this case, queries can only be answered for entities whose history has been maintained through chains. This requires that the database designer has previous knowledge of all possible queries that will involve version manipulation. Alternatively, these chains can be maintained for every entity and value in the database. Whereas the first alternative limits user exploratory activity, the second alternative requires a heavy overhead of pointers.

Supposing the historical chain of the designated area is available, then the queries are processed by the following procedure:

- (i) find the area and its vegetation cover in the present;
- (ii) follow back pointers of this area and cover, retrieving past information and rebuilding past states;
- (iii) produce the time series analysis desired.

Finally, the third query requires that not only actual historical chains be maintained, but also prevision chains for the same entity. This, again, complicates the housekeeping algorithms.

Pointer version mechanisms soon become too cumbersome to manage when each time period contains many entities that vary in different ways, as is typically the case in geographical applications.

It is interesting to compare these approaches to the DBV solution. The snapshot approach favors users who need to access entire contexts. However, it does not automatically support navigation through these contexts, since there is no sharing of entities across the snapshots. For instance, any update to an entity in a snapshot must be manually performed by the user in the other snapshots. The historical chain approach is geared towards management of versions of individual entities, through manipulation of their chains. However, it does not automatically support the building of contexts, which must be performed by the user.

The DBV model, on the other hand, automatically supports both working within a context and comparing entity versions across contexts. This is achieved thanks to the fact that it separates the logical versioning from the physical versioning level. It thus combines the advantages of the other approaches, without the inconvenients: it does not imply the waste of space of the snapshot approach, and neither does it demand the complex computation procedures of the historical chain approach.



## 6 Conclusions

This paper presented a solution for the management of evolution of georeferenced data in GIS which consists in using the DBV version mechanism. This solution allows the development of automated tools to keep track of different versions of the same georeferenced entity through time. This facility enables users to create and manage alternative scenarios, as well as to keep track of temporal data evolution. This type of support has so far been unavailable in commercial GIS, though required by different kinds of planning applications.

The use of a version mechanism, as pointed out in this paper, seems to be an obvious choice to cope with GIS users' demands. However, GIS databases do not consider this type of facility, since the handling of georeferenced data presents in itself many problems. Furthermore, available version management systems are complex and cannot readily satisfy GIS requirements.

The DBV mechanism, on the other hand, allows efficiently keeping track of data and schema versions in a database, with considerable savings in space and computation time, as compared to other database version mechanisms. It allows dissociating the issues of context and configuration consistency from version maintenance, which is not possible in other version models.

Its main advantage, from a GIS point of view, is that it allows users to access entire consistent database states for any given georeferenced entity. Thus users can create different scenarios by just modifying individual entities, and need not worry about keeping them within their appropriate context. This is achieved without additional overhead, by the appropriate management of version identifiers (as opposed to traditional mechanisms that require handling pointer chains). Finally, the DBV model is orthogonal to the underlying data model and to concurrency control, which are complicating factors in other version models.

### Acknowledgements

The research described in this paper was partially financed by grants FAPESP 91/2117-1, CNPq 453176/91, and CNPq 452357/93-4.

## References

- [Bat92] P. Batty. Exploiting relational database technology in a GIS. *Computers and Geosciences: An international journal*, 18(4):453–462, 1992.
- [BBA91] M. Borhani, J-P Barthès, and P. Anot. Versions in Object-Oriented Databases. Technical Report UTC/GI/DI/N 83, Université de Technologie de Compiègne, 1991.
- [CJ90] W. Cellary and G. Jomier. Consistency of Versions in Object-Oriented Databases. In *Proc. 16th VLDB*, pages 432–441, 1990.
- [CVJ94] W. Cellary, G. Vossen, and G. Jomier. Multiversion Object Constellations: a new Approach to Support a Designer's Database Work. *To appear, Engineering with Computers*, 1994.
- [Ege92] M. Egenhofer. Why not SQL! *International Journal of Geographical Information Systems*, 6(2):71–86, 1992.

- [FCF92] A. Frank, I. Campari, and U. Formentini, editors. *Theories and Methods of Spatio-Temporal Reasoning in Geographic Space*. Lecture Notes in Computer Science 639. Springer-Verlag, 1992.
- [Flo91] R. Flowerdew. *Geographical Information Systems - volume I*, chapter Spatial Data Integration, pages 375–387. John Wiley and Sons, 1991.
- [Gut92] J. Guttenberg. Towards a Behavioral Theory of Regionalization. In *Proc International Conference on GIS - From Space to Territory: Theories and Methods of Spatial Reasoning*, Springer Verlag Lecture Notes in Computer Science 639, pages 110–121, 1992.
- [HKGW93] N. Hachem, K. Qiu, M. Gennert, and M. Ward. Managing Derived Data in the GAEA Scientific DBMS. In *Proc 19th VLDB*, pages 1–12, 1993.
- [Kat90] R. H. Katz. Toward a Unified Framework for Version Modelling in Engineering Databases. *ACM Computing Surveys*, 22(4):375–408, 1990.
- [KS92] W. Kafer and H. Schoning. Mapping a Version Model to a Complex-Object Data Model. In *Proc IEEE Data Engineering Conference*, pages 348–357, 1992.
- [KSW86] P. Klahold, G. Schlageter, and W. Wilkes. A General Model for Version Management in Databases. In *Proc XII VLDB*, pages 319–327, 1986.
- [LST93] G. Landau, J. Schmidt, and V. Tsotras. Efficient Support of Historical Queries for Multiple Lines of Evolution. In *Proc. IEEE Data Engineering Conference*, pages 319–326, 1993.
- [MJ93] C. B. Medeiros and G. Jomier. Managing Alternatives and Data Evolution in GIS. In *Proc. ACM/ISCA Workshop on Advances in Geographic Information Systems*, pages 34–37, 1993.
- [MP94] C. B. Medeiros and F. Pires. Databases for GIS. *ACM Sigmod Record*, 23(1):107–115, 1994.
- [NTE92] R. Newell, D. Theriault, and M. Easterfieldy. Temporal GIS - modeling the evolution of spatial data in time. *Computers and Geosciences: An international journal*, 18(4):427–434, 1992.
- [Peu93] D. Peuquet. What, Where and When - a Conceptual Basis for Design of Spatiotemporal GIS Databases. In *Proc. ACM/ISCA Workshop on Advances in Geographic Information Systems*, pages 117–122, 1993.
- [PMB93] F. Pires, C. B. Medeiros, and A. Barros. Modelling Geographic Information Systems using an Object Oriented Framework. In *Proc XIII International Conference of the Chilean Computer Science Society*, pages 217–232, 1993.
- [Sam89] H. Samet. *The Design and Analysis of Spatial Data Structures*. Addison-Wesley, 1989.
- [Sno92] R. Snodgrass. Temporal Databases. In *Proc International Conference on GIS - From Space to Territory: Theories and Methods of Spatial Reasoning*, Springer Verlag Lecture Notes in Computer Science 639, pages 22–63, 1992.
- [TG92] V. Tsotras and B. Gopinath. Optimal Versioning of Objects. In *Proc IEEE Data Engineering Conference*, pages 358–365, 1992.