

# Modeling concept evolution: a historical perspective

Flavio Rizzolo, Yannis Velegarakis, John Mylopoulos, and Sjarhei Bykau

University of Trento, Trento, 38100, Italy  
{flavio, velgias, jm, bykau}@disi.unitn.eu

**Abstract.** The world is changing, and so must the data that describes its history. Not surprisingly, considerable research effort has been spent in Databases along this direction, covering topics such as temporal models and schema evolution. A topic that has not received much attention, however, is that of concept evolution. For example, Germany (instance-level concept) has evolved several times in the last century as it went through different governance structures, then split into two national entities that eventually joined again. Likewise, a caterpillar is transformed into a butterfly, while a mother becomes two (maternally-related) entities. As well, the concept of Whale (a class-level concept) changed over the past two centuries thanks to scientific discoveries that led to a better understanding of what the concept entails. In this work, we present a formal framework for modeling, querying and managing such evolution. In particular, we describe how to model the evolution of a concept, and how this modeling can be used to answer historical queries of the form “How has concept X evolved over period Y”. Our proposal extends an RDF-like model with temporal features and evolution operators. Then we provide a query language that exploits these extensions and supports historical queries.

## 1 Introduction

Conceptual modelling languages – including the ER Model, UML class diagrams and Description Logics – are all founded on a notion of “entity” that represents a “thing” in the application domain. Although the state of an entity can change over its lifetime, entities themselves are atomic and immutable. Unfortunately, this feature prevents existing modeling languages from capturing phenomena that involve the evolution of an entity into something else, such as a caterpillar becoming a butterfly, or Germany splitting off into two Germanies right after WWII. In these cases, there is general agreement that an entity evolves into one or more different entities. Moreover, there is a strong relationship between the two, which is not captured by merely deleting one entity and then creating another. Evolution applies not only at the instance level, but also at a class-level. For example, the concept of Whale evolved over the past two centuries in the sense that whales were considered some sort of fish, but are now recognized as mammals. We are interested in modeling this kind of evolution relationship, both at the instance- and at the class-level.

In Databases, considerable amount of research effort has been spent on the development of models, techniques and tools for modeling and managing data changes. These range from data manipulation languages, and maintenance of views under changes [1],

to schema evolution [2] and mapping adaptation [3]. To cope with the history of data changes, temporal models have been proposed for the relational [4] and ER [5] models, for semi-structured data [6], XML [7] and for RDF [8]. Almost in its entirety, existing work on data changes is based on a data-oriented point of view. It aims at recording and managing changes that are taking place at the values of the data. What has been completely overlooked are other types of changes, such as an entity evolving/mutating into another, or an entity “splitting off” into several others.

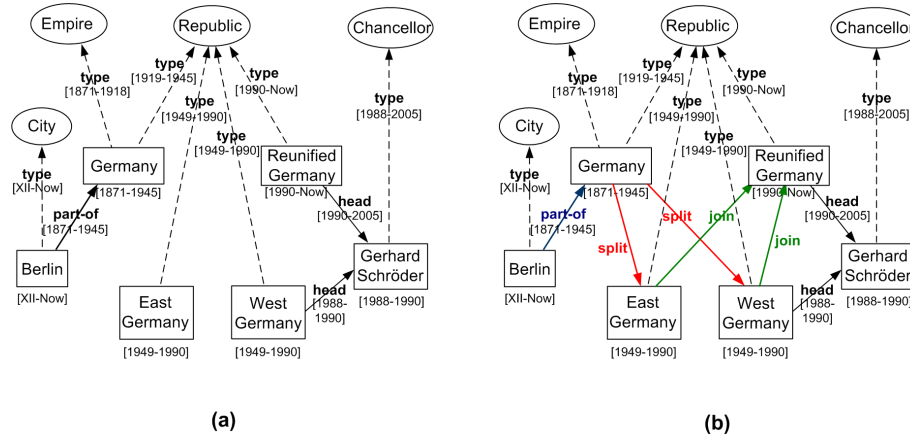
In this work, we propose a framework for modeling the evolution of an entity into one or more different entities, as well as the inter-dependencies among the pre- and post-evolution entities. The framework allows posing new kinds of queries that previously could not have been expressed. For instance, we aim at supporting queries of the form: How has an entity evolved over time? From what other entities has it evolved and into what others has it resulted? What other entities have affected its evolution over time? What concepts are indirectly related to it and how? These kinds of queries are of major importance for many interesting areas:

**Historical Studies.** Modern historians are interested in studying the history of human achievements, events and important persons. In addition, they want to understand how systems, tools, concepts, and techniques have evolved throughout the history. For them it is not enough to query a data source for a specific moment in history. They need to ask questions on how concepts and the relationships that exist between them have changed over time. Historians may be interested in the evolution of countries like Germany, with respect to territory, political division, etc. Alternatively, they may want to study the evolution of scientific topics, e.g., how the concept of biotechnology has evolved from its beginnings as an agricultural technology to the current notion that is coupled to genetics and molecular biology.

**Entity Management.** Web application and integration systems are progressively moving from tuple and value-based towards entity-based solutions, i.e., systems in which the basic data unit is an entity, independently of its logical modeling [9]. Furthermore, web integration systems, in order to achieve interoperability, may need to provide unique identification for the entities in the data they exchange [10]. Unfortunately, entities do not remain static over time. They evolve, merge, split, get created and disappear. Knowing the history of each entity, i.e., how it has been formed and from what, paves the ground for successful entity management solutions and affective information exchange.

**Life Sciences** One of the Biology fields is the study of the evolution of the species since life started on earth. To better understand the secrets of nature, it is important to model how the different species have evolved, from what, if and how they have disappeared, and when.

In this work we present a framework for modeling the evolution of concepts over time and the evolving relationships among them. In particular, our contributions are the following: (i) we consider a conceptual model enhanced with the concept of a lifetime of a class/individual and even of a relationship, as first class citizen, (ii) we further extend the temporal model [8] with consistency conditions and additional constructs to model merges, splits, and other forms of evolution among concepts and individuals; (iii) we introduce a query language that allows the answering of queries regarding the lifetime of concepts as well as the way they have evolved over time, along with other



**Fig. 1.** The concepts of Germany, in a temporal model (a) and in our evolutionary model (b)

associated (via evolution) concepts; and finally, (iv) we present a case study in which we have applied our framework and we describe our findings.

## 2 A Motivating Example

Consider the knowledge base of a historian that records information about countries and their political governance. A fraction of that information modeling a part of the history of Germany is illustrated in Figure 1. In it, the status of Germany at different times in history has been modeled through different individuals or through different instantiations. In particular, since 1871 and until 1945, Germany has been an empire and later a republic. This change is modeled by the multiple instantiation of Germany to Empire and Republic respectively. Shortly after the end of the World War II, Germany was split in four zones<sup>1</sup> that in 1949 formed the East and West Germany. These two parts lasted until 1990, when they were merged to form the republic of Germany as we know it today, which is modeled through the individual Reunified Germany.

To model the validity of each state of Germany at the different periods, a temporal model similar to temporal RDF [8] can be used. The model associates to each concept or individual a specific time frame. The time frames assigned to the individuals that model Germany are illustrated in Figure 1 through the intervals next to each individual. The same applies to every property, instance and subclass relationship. Note, for example, how the instantiation relationship of Germany to Empire has a temporal interval from 1871 to 1918, while the one to Republic has a temporal interval from 1918 to 1945. It is important for such a knowledge base to contain no inconsistencies, i.e., situations like having an instantiation relationship with a temporal interval bigger than the interval of the class it instantiates. Although temporal RDF lacks this kind of axiomatization, a mechanism for consistency checking needs to be in place for the accurate representation

<sup>1</sup> The information of the four zones is not illustrated in the Figure 1

of concepts and individuals in history. Our solution provides an axiomatization of the temporal RDF that guarantees the consistency of the historical information recorded in a knowledge base.

Consider now the case of a historian that is interested in studying the history of Germany. Using traditional query mechanisms, the historian will only be able to retrieve the individual Germany. Using keyword based techniques, she may be able to also retrieve the remaining individuals modeling Germany at different times in history, but this under the assumption that each such individual contains the keyword “Germany”. Applying terminology evolution [11], it is possible to infer that the four terms for Germany, i.e., Germany, East Germany, West Germany, and Reunified Germany refer to related concepts regardless of the keywords that appear in the terms. Yet, in neither case the historian will be able to reconstruct how the constituents of Germany have changed over time. She will not be able to find that the East and West Germany were made by splitting the pre-war Germany and its parts, neither that East and West Germany were the same two that merged to form the Reunified Germany.

We propose the use of explicit constructs to allow the modeling of the sequential conceptual evolution in knowledge bases. A small illustration of such constructs (split, merge, etc.) can be found in Figure 1(b).

Assuming such constructs are in place, let the historian pose the question whether Berlin is part of the Reunified Germany. Since Berlin is modeled to be part of Germany and not of Reunified Germany, the historian will get a negative answer. Since the Reunified Germany was formed (through some intermediate steps) from the initial pre-war Germany, one would expect the answer to the query of whether Berlin is part of Reunified Germany to be positive. Furthermore, a typical historian query is the one asking for the list of all the leaders of a country in history. Gerhard Schroeder was a leader of the West Germany and of the Reunified Germany, but not of the pre-War Germany. However, the historian should be able to express queries that can collect the leaders of Germany from all its phases through time.

### 3 Temporal Knowledge Bases

We consider an RDF-like data model. The model is expressive enough to represent ER models and the majority of ontologies and schemas that are met in practice [12]. It does not include certain OWL Lite features such as *sameAs* or *equivalentClass*, since these features have been considered to be out of the main scope of this work and their omission does not restrict the functionality of the model.

We assume the existence of an infinite set of *resources*  $\mathcal{U}$ , each with a *unique resource identifier* (URIs), and a set of *literals*  $\mathcal{L}$ . A *property* is a relationship between two resources. Properties are considered resources.

We consider the existence of the special properties: `rdfs:type`, `rdfs:domain`, `rdfs:range`, `rdfs:subClassOf` and `rdfs:subPropertyOf`, which we denote for simplicity as `type`, `dom`, `rng`, `subc`, and `subp`, respectively. The set  $\mathcal{U}$  contains three special resources: `rdfs:Property`, `rdfs:Class` and `rdf:Thing`, which we denote for simplicity as `Prop`, `Class` and `Thing`, respectively. The semantics of these resources as well as the semantics of the special properties are those defined in RDFS [13].

Resources are described by a set of triples that form a knowledge base.

**Definition 1.** A knowledge base  $\Sigma$  is a tuple  $\langle U, L, T \rangle$ , where  $U \subseteq \mathcal{U}$ ,  $L \subseteq \mathcal{L}$ ,  $T \subseteq \mathcal{U} \times \mathcal{U} \times \{\mathcal{U} \cup \mathcal{L}\}$ , and  $U$  contains the resources `rdfs:Property`, `rdfs:Class`, and `rdf:Thing`. The set of classes of the knowledge base  $\Sigma$  is the set  $C = \{x \mid \exists \langle x, \text{type}, \text{rdfs:Class} \rangle \in T\}$ . Similarly, the set of properties is the set  $P = \{x \mid \exists \langle x, \text{type}, \text{rdfs:Property} \rangle \in T\}$ . The set  $P$  must contain the RDFS properties `type`, `dom`, `rng`, `subc`, and `subp`. A resource  $i$  is said to be an instance of a class  $c \in C$  (or of type  $c$ ) if  $\exists \langle i, \text{type}, c \rangle \in T$ . The set of instances is the set  $I = \{i \mid \exists \langle i, \text{type}, y \rangle \in T\}$ .

A knowledge base can be equivalently represented as a hypergraph, which we call *RDF graph*. Thus, the rest of the paper, we will use the terms knowledge base and RDF graph, equivalently.

**Definition 2.** An RDF graph of a knowledge base  $\Sigma$  is an hypergraph in which nodes represent resources and literals and the edges represent triples.

*Example 1.* Figure 1(a) is an illustration of an RDF Graph. The nodes Berlin and Germany represent resources. The edge labeled part-of between them represents the triple  $\langle \text{Berlin}, \text{part-of}, \text{Germany} \rangle$ . The label of the edge, i.e., part-of, represents a property.

To support the temporal dimension in our model, we adopt the approach of Temporal RDF [8] which extends RDF by associating to each triple a time frame. Unfortunately, this extension is not enough for our goals. We need to add time semantics not only to relationships between resources (what the triples represent), but also to resources themselves by providing temporal-varying classes and individuals. This addition and the consistency conditions we introduce below are our temporal extensions to the temporal RDF data model.

We consider time as a discrete, total order domain  $\mathbb{T}$  in which we define different granularities. Following [14], a *granularity* is a mapping from integers to *granules* (i.e., subsets of the time domain  $\mathbb{T}$ ) such that contiguous integers are mapped to non-empty granules and granules within one granularity are totally ordered and do not overlap. Days and months are examples of two different granularities, in which each granule is a specific day in the former and a month in the latter. Granularities define a lattice in which granules in some granularities can be aggregated in larger granules in *coarser* granularities. For instance, months are a coarser granularity than days because every granule in the former (a month) is composed of an integer number of granules in the latter (days). In contrast, months are not coarser (nor finer) than weeks because not every month is composed by an integer number of weeks (and obviously no week is compose of an integer number of months).

Even though we model time as a point-based temporal domain, we use intervals as abbreviations of sets of instants whenever possible. An ordered pair  $[a, b]$  of time points, with  $a, b$  granules in a granularity, and  $a \leq b$ , denotes the closed interval from  $a$  to  $b$ . As in most temporal models, the current time point will be represented with the distinguished word *Now*. We will use the symbol  $\mathcal{T}$  to represent the infinite set of all the possible temporal intervals over the temporal domain  $\mathbb{T}$ , and the expressions *i.start* and *i.end* to refer to the starting and ending time points of an interval  $i$ . Given two

intervals  $i_1$  and  $i_2$ , we will denote by  $i_1 \sqsubseteq i_2$  the containment relationship between the intervals in which  $i_2.start \leq i_1.start$  and  $i_1.end \leq i_2.end$ .

**Definition 3.** A temporal knowledge base  $\Sigma_T$  is a tuple  $\langle U, L, T, \tau \rangle$ , where  $\langle U, L, T \rangle$  is a knowledge base and  $\tau$  is function that maps every resource  $r \in U$  to a temporal interval in  $T$ . The temporal interval is also referred to as the lifespan of the resource. The expressions  $r.start$  and  $r.end$  denote the start and end points and of the interval, respectively. The temporal graph of  $\Sigma_T$  is the RDF graph of  $\langle U, L, T \rangle$  enhanced with the temporal intervals on the edges and nodes.

For a temporal knowledge base to be semantically meaningful, the lifespan of the resources need to satisfy certain conditions. For instance, it is not logical to have an individual with a lifespan that does not contain any common time points with the lifespan of the class it belongs to. Temporal RDF does not provide such a mechanism, thus, we are introducing the notion of a consistent temporal knowledge base.

**Definition 4.** A consistent temporal knowledge base is a temporal knowledge base  $\Sigma_\tau = \langle U, L, T, \tau \rangle$  that satisfies the following conditions:

1.  $\forall r \in L \cup \{\text{Prop, Class, Thing, type, dom, rng, subc, subp}\}: \tau(r) = [0, \text{Now}]$ ;
2.  $\forall \langle d, p, r \rangle \in T : \tau(\langle d, p, r \rangle) \sqsubseteq \tau(d)$  and  $\tau(\langle d, p, r \rangle) \sqsubseteq \tau(r)$ ;
3.  $\forall \langle d, p, r \rangle \in T$  with  $p \in \{\text{type, subc, subp}\}: \tau(d) \sqsubseteq \tau(r)$ .

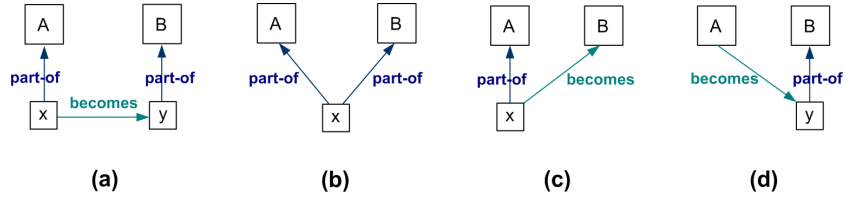
Intuitively, literals and the special resources and properties defined in RDFS need to be valid during the entire lifespan of the temporal knowledge base, which is  $[0, \text{Now}]$  (Condition 1). In addition, the lifespan of a triple needs to be within the lifespan of the resources that the triple associates (Condition 2). Finally, the lifespan of a resource has to be within the lifespan of the class the resource instantiates, and any class or property needs to be within the lifespan of its superclasses or superproperties (Condition 3).

## 4 Modeling evolution

Apart from the temporal dimension that was previously described, two new dimensions need to be taken into consideration to successfully model evolution: the mereological and the causal.

Mereology [15] is a sub-discipline in philosophy that deals with the ontological investigation of the part-whole relationships. It is used in our model to capture the parthood relationship between concepts in a way that is carried forward as concepts evolve. Such a relationship is modeled through the introduction of the special property **part-of**, which is reflexive, antisymmetric and transitive. A property **part-of** is defined from a resource  $x$  to a resource  $y$  if the concept modeled by resource  $x$  is part of the concept modeled by resource  $y$ . Apart from this special semantics, **part-of** behaves as any other property in a temporal knowledge base. For readability and presentation reasons, we may use the notation  $x \xrightarrow{\text{part-of}} y$  to represent the existence of a triple  $\langle x, \text{part-of}, y \rangle$  in the set  $T$  of a temporal knowledge base  $\Sigma_\tau$ .

To capture the causal relationships, i.e., the interdependency between two resources whose lifespans do not overlap, we additionally introduce the notion of **becomes**,



**Fig. 2.** Liaison examples

which is an antisymmetric and transitive relation. For similar reasons as before, we may use the notation  $x \xrightarrow{\text{becomes}} y$  to represent the fact that  $(x, y) \in \mathbf{becomes}$ . Intuitively,  $x \xrightarrow{\text{becomes}} y$  means that the concept modeled by resource  $y$  originates from the concept modeled by resource  $x$ . In addition to not overlapping, we require that  $\tau(x).end < \tau(y).starts$ .

To effectively model evolution, we introduce the notion of a *liaison*. A *liaison* between two concepts is another concept that keeps the former two linked together in time by means of **part-of** and **becomes**. A liaison is part of at least one of the concepts it relates and has some causal relationship to a part of the other. This situation is graphically depicted in Figure 2 (a). The boxes  $A$  and  $B$  represent the two main concepts whereas the  $x$  and  $y$  represent two of their respective parts. Figure 2 (b) illustrates the special case in which  $x$  and  $y$  are actually the same concept. Figure 2 (c) (respectively, (d)) shows the special case in which  $y$  (respectively,  $x$ ) is exactly the whole of  $B$  (respectively,  $A$ ) rather than a part of it.

**Definition 5 (Liaison).** Let  $A, B$  be two concepts of a temporal knowledge base with  $\tau(A).start < \tau(B).start$ , and  $x, y$  concepts for which  $x \xrightarrow{\text{part-of}} A$  and  $y \xrightarrow{\text{part-of}} B$ . A concept  $x$  (or  $y$ ) is said to be a liaison between  $A$  and  $B$  if either  $x \xrightarrow{\text{becomes}} y$ , or  $x \xrightarrow{\text{part-of}} B$ , or  $x \xrightarrow{\text{becomes}} B$ , or  $A \xrightarrow{\text{becomes}} y$ .

To model the different kinds of evolution events that may exist, we introduce four evolution terms: *join*, *split*, *merge*, and *detach*.

**[join]** The *join* term, denoted as  $\mathbf{join}(c_1 \dots c_n, c, t)$ , models the fact that when the new concept  $c$  is born at time  $t$ , every part of it comes from a part (or the whole) of some concept in  $\{c_1, \dots, c_n\}$ . In particular:

- $\tau(c).start = t$
- $\forall x$  s.t.  $x \xrightarrow{\text{part-of}} c: \exists c_i$  s.t.  $x$  is a liaison between  $c_i$  and  $c$ , with  $1 \leq i \leq n$

**[split]** The *split* term, denoted as  $\mathbf{split}(c, c_1 \dots c_n, t)$ , is used to model the fact that when the lifespan of a concept  $c$  ends at a time  $t$ , every part of  $c$  becomes either one of the new independent concepts  $c_1, \dots, c_n$ , or one of their parts. In particular:

- $\tau(c).end = t$
- $\forall x$  s.t.  $x \xrightarrow{\text{part-of}} c: \exists c_i$  s.t.  $x$  is a liaison between  $c$  and  $c_i$ , with  $1 \leq i \leq n$

**[merge]** The *merge* term, denoted as  $\mathbf{merge}(c, c', t)$ , is used to model the fact that the lifespan of a concept  $c$  ends at a time  $t$  and at least a part of it becomes part of an existing concept  $c'$ . In particular:

- $\tau(c).end=t$
- $\exists x$  s.t.  $x \xrightarrow{part.of} c$ :  $x$  is a liaison between  $c$  and  $c'$

**[detach]** The *detach* term, denoted as  $\mathbf{detach}(c, c', t)$ , is used to model the fact that the new concept  $c'$  is formed at a time  $t$  with at least one part from  $c$ . In particular:

- $\tau(c').starts=t$
- $\exists x$  s.t.  $x \xrightarrow{part.of} c$ :  $x$  is a liaison between  $c$  and  $c'$

We record the **becomes** relation and the evolution terms in the temporal knowledge base as *evolution triples*  $\langle c, term, c' \rangle$ , where *term* is one of the special evolution properties **becomes**, **join**, **split**, **merge**, and **detach**. Evolution properties are meta-temporal, i.e., they describe how the temporal model changes, and thus their triples do not need to satisfy the consistency conditions in Definition 4. A temporal knowledge base with a set of evolution properties and triples defines an *evolution base*.

**Definition 6.** An evolution base  $\Sigma_T^E$  is a tuple  $\langle U, L, T, E, \tau \rangle$ , where  $\langle U, L, T, \tau \rangle$  is a temporal knowledge base,  $U$  contains a set of evolution properties, and  $E$  is a set of evolution triples. The evolution graph of  $\Sigma_T^E$  is the temporal graph of  $\langle U, L, T, \tau \rangle$  enhanced with edges representing the evolutions triples.

The time in which the evolution event took place does not need to be recorded explicitly in the triple since it can be retrieved from the lifespan of the involved concepts. For instance,  $\mathbf{detach}(\text{Kingdom of the Netherlands}, \text{Belgium}, 1831)$  is modeled as the triple:  $\langle \text{Kingdom of the Netherlands}, \mathbf{detach}, \text{Belgium} \rangle$  with  $\tau(\text{Belgium}).start = 1831$ .

For recording evolution terms that involve more than two concepts, e.g. the **join**, multiple triples are needed. We assume that the terms are indexed by their time, thus, the set of (independent) triples that belong to the same terms can be easily detected since they will all share the same start or end time in the lifespan of the respective concept. For instance,  $\mathbf{split}(\text{Germany}, \text{East Germany}, \text{West Germany}, 1949)$  is represented in our model through the triples  $\langle \text{Germany}, \mathbf{split}, \text{East Germany} \rangle$  and  $\langle \text{Germany}, \mathbf{split}, \text{West Germany} \rangle$  with  $\tau(\text{East Germany}).start = \tau(\text{West Germany}).start = 1949$ .

Note that the evolution terms may entail facts that are not explicitly represented in the knowledge base. For instance, the split of Germany into West and East implies the fact that Berlin, which is explicitly defined as part of Germany, becomes part of either East or West. This kind of reasoning is beyond the scope of the current work.

## 5 Query Language

We define a navigational query language to traverse temporal and evolution edges in an evolution graph. This language is analogous to nSPARQL [16], a language that



$$\begin{aligned}
\mathcal{E}[\mathbf{self}] &:= \{\langle x, x, \tau(x) \rangle \mid x \in U \cup L\} \\
\mathcal{E}[\mathbf{self}::r] &:= \{\langle r, r, \tau(r) \rangle\} \\
\mathcal{E}[\mathbf{next}] &:= \{\langle x, y, \tau(t) \rangle \mid t = \langle x, z, y \rangle \in T\} \\
\mathcal{E}[\mathbf{next}::r] &:= \{\langle x, y, \tau(t) \rangle \mid t = \langle x, r, y \rangle \in T\} \\
\mathcal{E}[\mathbf{edge}] &:= \{\langle x, y, \tau(t) \rangle \mid t = \langle x, y, z \rangle \in T\} \\
\mathcal{E}[\mathbf{edge}::r] &:= \{\langle x, y, \tau(t) \rangle \mid t = \langle x, y, r \rangle \in T\} \\
\mathcal{E}[\mathbf{node}] &:= \{\langle x, y, \tau(t) \rangle \mid t = \langle z, x, y \rangle \in T\} \\
\mathcal{E}[\mathbf{node}::r] &:= \{\langle x, y, \tau(t) \rangle \mid t = \langle r, x, y \rangle \in T\} \\
\mathcal{E}[\mathbf{e-edge}] &:= \{\langle x, \mathbf{e-axis}, [0, Now] \rangle \mid t = \langle x, \mathbf{e-axis}, z \rangle \in E\} \\
\mathcal{E}[\mathbf{e-node}] &:= \{\langle \mathbf{e-axis}, y, \tau(t) \rangle \mid t = \langle z, \mathbf{e-axis}, y \rangle \in E\} \\
\mathcal{E}[\mathbf{e-axis}] &:= \{\langle x, y, [0, Now] \rangle \mid \exists t = \langle x, \mathbf{e-axis}, y \rangle \in E\} \\
\mathcal{E}[\mathbf{forward}] &:= \bigcup_{\mathbf{e-axis}} \mathcal{E}[\mathbf{e-axis}] \\
\mathcal{E}[\mathbf{backward}] &:= \bigcup_{\mathbf{e-axis}} \mathcal{E}[\mathbf{e-axis}^{-1}] \\
\mathcal{E}[\mathbf{self}::[exp]] &:= \{\langle x, x, \tau(x) \cap I \rangle \mid x \in U \cup L, \exists \langle x, z, I \rangle \in \mathcal{P}[[exp]], \tau(x) \cap I \neq \emptyset\} \\
\mathcal{E}[\mathbf{next}::[exp]] &:= \{\langle x, y, \tau(t) \cap I \rangle \mid t = \langle x, z, y \rangle \in T, \exists \langle z, w, I \rangle \in \mathcal{P}[[exp]], \tau(t) \cap I \neq \emptyset\} \\
\mathcal{E}[\mathbf{edge}::[exp]] &:= \{\langle x, y, \tau(t) \cap I \rangle \mid t = \langle x, y, z \rangle \in T, \exists \langle z, w, I \rangle \in \mathcal{P}[[exp]], \tau(t) \cap I \neq \emptyset\} \\
\mathcal{E}[\mathbf{node}::[exp]] &:= \{\langle x, y, \tau(t) \cap I \rangle \mid t = \langle z, x, y \rangle \in T, \exists \langle z, w, I \rangle \in \mathcal{P}[[exp]], \tau(t) \cap I \neq \emptyset\} \\
\mathcal{E}[\mathbf{axis}^{-1}] &:= \{\langle x, y, \tau(t) \rangle \mid \langle y, x, \tau(t) \rangle \in \mathcal{E}[\mathbf{axis}]\} \\
\mathcal{E}[\mathbf{t-axis}^{-1}::r] &:= \{\langle x, y, \tau(t) \rangle \mid \langle y, x, \tau(t) \rangle \in \mathcal{E}[\mathbf{t-axis}::r]\} \\
\mathcal{E}[\mathbf{t-axis}^{-1}::[exp]] &:= \{\langle x, y, \tau(t) \rangle \mid \langle y, x, \tau(t) \rangle \in \mathcal{E}[\mathbf{t-axis}::[exp]]\} \\
\mathcal{E}[\mathbf{exp}[I]] &:= \{\langle x, y, I \cap I' \rangle \mid \langle x, y, I' \rangle \in \mathcal{E}[[exp]] \text{ and } I \cap I' \neq \emptyset\} \\
\mathcal{E}[\mathbf{exp}/\mathbf{e-exp}] &:= \{\langle x, y, I_2 \rangle \mid \exists \langle x, z, I_1 \rangle \in \mathcal{E}[[exp]], \exists \langle z, y, I_2 \rangle \in \mathcal{E}[[\mathbf{e-exp}]]\} \\
\mathcal{E}[\mathbf{exp}/\mathbf{t-exp}] &:= \{\langle x, y, I_1 \cap I_2 \rangle \mid \exists \langle x, z, I_1 \rangle \in \mathcal{E}[[exp]], \exists \langle z, y, I_2 \rangle \in \mathcal{E}[[\mathbf{t-exp}]] \text{ and } \\
&\quad I_1 \cap I_2 \neq \emptyset\} \\
\mathcal{E}[\mathbf{exp}_1|\mathbf{exp}_2] &:= \mathcal{E}[[\mathbf{exp}_1] \text{ cup } \mathcal{E}[[\mathbf{exp}_2]] \\
\mathcal{E}[\mathbf{exp}^*] &:= \mathcal{E}[\mathbf{self}] \cup \mathcal{E}[[exp]] \cup \mathcal{E}[[exp/exp]] \cup \mathcal{E}[[exp/exp/exp]] \cup \dots \\
\mathcal{P}[\mathbf{e-exp}] &:= \mathcal{E}[[\mathbf{e-exp}]] \\
\mathcal{P}[\mathbf{t-exp}] &:= \mathcal{E}[[\mathbf{t-exp}]] \\
\mathcal{P}[\mathbf{t-exp}/\mathbf{exp}] &:= \{\langle x, y, I_1 \cap I_2 \rangle \mid \exists \langle x, z, I_1 \rangle \in \mathcal{E}[[\mathbf{t-exp}]], \exists \langle z, y, I_2 \rangle \in \mathcal{E}[[exp]] \text{ and } \\
&\quad I_1 \cap I_2 \neq \emptyset\} \\
\mathcal{P}[\mathbf{e-exp}/\mathbf{exp}] &:= \{\langle x, y, I_1 \rangle \mid \exists \langle x, z, I_1 \rangle \in \mathcal{E}[[\mathbf{e-exp}]], \exists \langle z, y, I_2 \rangle \in \mathcal{E}[[exp]]\} \\
\mathcal{P}[\mathbf{exp}_1|\mathbf{exp}_2] &:= \mathcal{E}[[\mathbf{exp}_1|\mathbf{exp}_2]] \\
\mathcal{P}[\mathbf{exp}^*] &:= \mathcal{E}[[\mathbf{exp}^*]] \\
\mathbf{t-exp} &\in \{\mathbf{t-axis}, \mathbf{t-axis}::r, \mathbf{t-axis}::[exp], \mathbf{t-axis}[I]\} \text{ and} \\
\mathbf{e-exp} &\in \{\mathbf{e-axis}, \mathbf{e-axis}::[exp], \mathbf{e-axis}[I], \mathbf{forward}, \mathbf{backward}\}
\end{aligned}$$

**Fig. 3.** Formal semantics of nested evolution expressions

extends SPARQL with navigational capabilities based on nested regular expressions. nSPARQL uses four different axes, namely **self**, **next**, **edge**, and **node**, for navigation on an RDF graph and node label testing. We have extended the nested regular expressions constructs of nSPARQL with temporal semantics and a set of five *evolution axes*, namely **join**, **split**, **merge**, **detach**, and **becomes** that extend the traversing capabilities of nSPARQL to the evolution edges. The language is defined according to

the following grammar:

$$exp := \text{axis} \mid \text{t-axis} :: a \mid \text{t-axis} :: [exp] \mid exp[I] \mid exp/exp \mid exp|exp \mid exp^*$$

where  $a$  is a node in the graph,  $I$  is a time interval, and **axis** can be either **forward**, **backward**, **e-edge**, **e-node**, a **t-axis** or an **e-axis**, with **t-axis**  $\in$  {**self**, **next**, **edge**, **node**} and **e-axis**  $\in$  {**join**, **split**, **merge**, **detach**, **becomes**}.

The evaluation of an evolution expression  $exp$  is given by the semantic function  $\mathcal{E}$  defined in Figure 3.  $\mathcal{E}[[exp]]$  returns a set of tuples of the form  $\langle x, y, I \rangle$  such that there is a path from  $x$  to  $y$  satisfying  $exp$  during interval  $I$ . For instance, in the evolution base of Figure 1,  $\mathcal{E}[[\text{self} :: \text{Germany}/\text{next} :: \text{head}/\text{next} :: \text{type}]]$  returns the tuple  $\langle \text{Germany}, \text{Chancellor}, [1988, 2005] \rangle$ . It is also possible to navigate an edge from a node using the edge axis and to have a nested expression  $[exp]$  that functions as a predicate that the preceding expression must satisfy. For example,  $\mathcal{E}[[\text{self}[\text{next} :: \text{head}/\text{self} :: \text{Gerhard Schröder}]]]$  returns  $\langle \text{Reunified Germany}, \text{Reunified Germany}, [1990, 2005] \rangle$  and  $\langle \text{West Germany}, \text{West Germany}, [1988, 1990] \rangle$ .

In order to support evolution expressions, we need to extend nSPARQL triple patterns with temporal and evolution semantics. In particular, we redefine the evaluation of an nSPARQL triple pattern  $(?X, exp, ?Y)$  to be the set of triples  $\langle x, y, I \rangle$  that result from the evaluation of the evolution expression  $exp$ , with the variables  $X$  and  $Y$  bound to  $x$  and  $y$ , respectively. In particular:

$$\mathcal{E}[[(?X, exp, ?Y)]] := \{(\theta(?X), \theta(?Y)) \mid \theta(?X) = x \text{ and } \theta(?Y) = y \text{ and } \langle x, y, I \rangle \in \mathcal{E}[[exp]]\}$$

Our language includes all nSPARQL operators such as AND, OPT, UNION and FILTER with the same semantics as in nSPARQL. For instance:

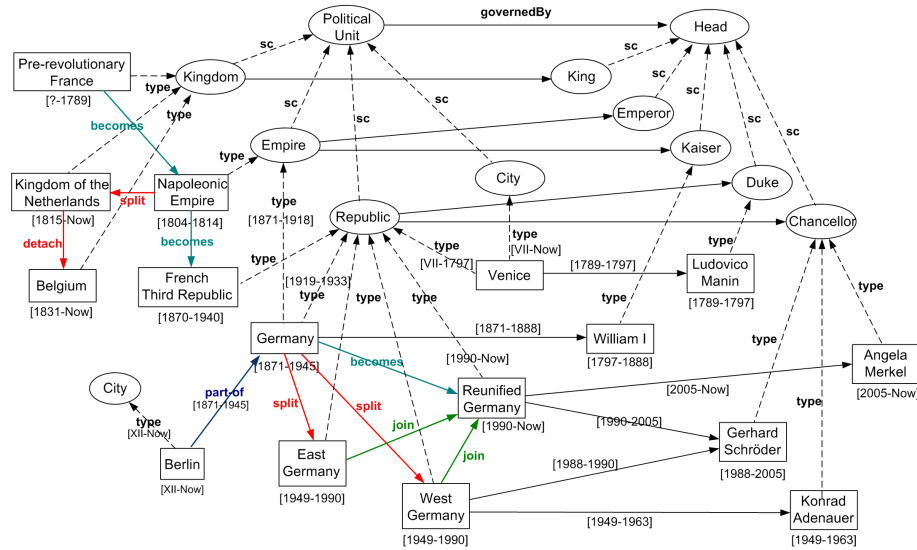
$$\mathcal{E}[[P_1 \text{ AND } P_2]] := \mathcal{E}[[P_1]] \bowtie \mathcal{E}[[P_2]]$$

where  $P_1$  and  $P_2$  are triple patterns and  $\bowtie$  is the join on the variables  $P_1$  and  $P_2$  have in common. A complete list of all the nSPARQL operators and their semantics can be found in [16].

## 6 Application scenarios

Consider an evolution base that models how countries have changed over time in terms of territory, political division, type of government, etc. Classes are represented by ovals and instances by boxes. A small fragment of that evolution base is illustrated as a graph in Figure 4.

Germany is a concept that has changed several times along history. The country was unified as a nation-state in 1871 and the concept of Germany first appears in our historical knowledge base as Germany at instant 1871. After WWII, Germany was divided into four military zones (not shown in the figure) that were merged into West and East Germany in 1949. This is represented with two split edges from the concept of Germany to the concepts of West Germany and East Germany. The country was finally reunified in 1990, which is represented by the coalescence of the West Germany and East Germany concepts into Unified Germany via two merge edges. These merge and



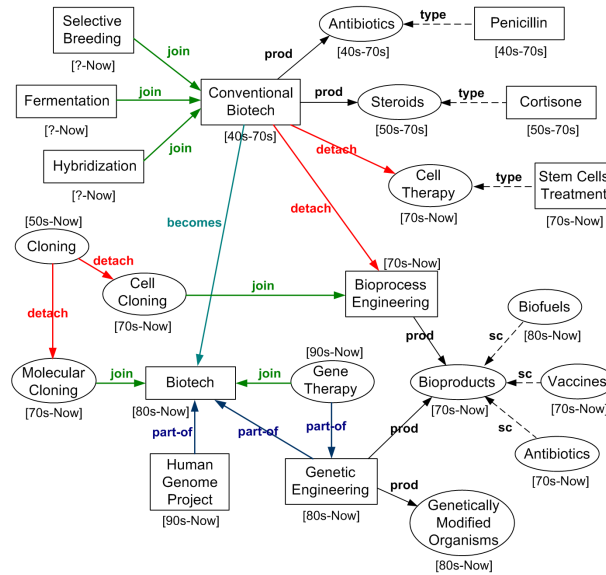
**Fig. 4.** The evolution of the concepts of Germany and France (full black lines represent governedBy properties)

split constructs are also defined in terms of the parts of the concepts they relate. For instance, a part-of property indicates that Berlin was part of Germany during [1871, 1945]. Since that concept of Germany existed until 1945 whereas Berlin exists until today, the part-of relation is carried forward by the semantics of split and merge into the concept of Reunified Germany. Consider now a historian who is interested in finding answers to a number of evolution-related queries.

**[Example Query 1]:** How has the notion of Germany changed over the last two centuries in terms of its constituents, government, etc.? The query can be expressed in our extended query language as follows:

```
Select ?Y, ?Z, ?W
(?X, self :: Reunified Germany / backward * [1800, 2000] /, ?Y) AND
(?Y, edge, ?Z) AND (?Z, edge, ?W)
```

The query first binds  $?X$  to Reunified Germany and then follows all possible evolution axes backwards in the period [1800, 2000]. All concepts bound to  $?Y$  are in an evolution path to Reunified Germany, namely Germany, West Germany, and East Germany. Note that, since the semantics of an  $*$  expression includes **self** (see Figure 3), then Reunified Germany will also bind  $?Y$ . The second triple returns in  $?Z$  the name of the properties of which  $?Y$  is the subject, and finally the last triple returns in  $?W$  the objects of those properties. By selecting  $?Y, ?Z, ?W$  in head of the query, we get all evolutions of Germany together with their properties.



**Fig. 5.** The evolution of the concept of Biotechnology

**[Example Query 2]:** Who was the head of the German government before and after the unification of 1990? The query can be expressed as follows:

```
Select ?Y
(?X, self :: Reunified Germany / join-1[1990] / next :: head[1990], ?Y) AND
(?Z, self :: Reunified Germany / next :: head[1990], ?Y)
```

The first triple finds all the heads of state of the Reunified Germany before the unification by following  $\text{join}^{-1}[1990]$  and then following  $\text{next} :: \text{head}[1990]$ . The second triple finds the heads of state of the Reunified Germany. Finally, the join on  $?Y$  will bind the variable only to those heads of state that are the same in both triples, hence returning the one before and after the mentioned unification.

Consider now the evolution of the concept of biotechnology from a historical point of view. According to historians, biotechnology got its current meaning (related to molecular biology) only after the 70s. Before that, the term biotechnology was used in areas as diverse as agriculture, microbiology, and enzyme-based fermentation. Even though the term “biotechnology” was coined in 1919 by Karl Ereky, a Hungarian engineer, the earliest mentions of biotechnology in the news and specialized media refer to a set of ancient techniques like selective breeding, fermentation and hybridization. Even though from the 70s the dominant meaning of biotechnology has been closely related to genetics, it is possible to find a large body of news and other media articles from the 60s to the 80s that use the term biotechnology to refer to an environmentally friendly technological orientation entirely unrelated to genetics but closely related to bioprocess engineering. Not only the use of the term changed from the 60s to the 90s, but also the two different meanings coexisted in the media for almost two decades.

Figure 5 illustrates the evolution of the notion of biotechnology since the 40s. The notions of Selective breeding, Fermentation and Hybridization existed from an indeterminate time until now and in the 40s joined the new topic of Conventional Biotech, which groups ancient techniques like the ones mentioned above. Over the next decades, Conventional Biotech started to include more modern therapies and products such as Cell Therapy, Penicillin and Cortisone. At some point in the 70s, the notions of Cell Therapy and Bioprocess Engineering matured and detached from Conventional Biotech becoming independent concepts. The three concepts coexisted in time during part of the 70s and the 80s. During the 80s, the notion of Conventional Biotech stopped being used and all its related concepts became independent topics. In parallel to this, the new topic of Biotech started to take shape. We could see Biotech as an evolution of the former Conventional Biotech but using Genetic Engineering instead of conventional techniques. Concepts and terms related to the Biotech and Genetic Engineering topics are modeled with a part-of property.

**[Example Query 3]:** Is the academic discipline of biotechnology a wholly new technology branch or has it derived from the combination of other disciplines? Which ones and how? The query requires to follow evolution paths and return the traversed triples in addition to the nodes in order to answer the question of “how”. The query is expressed in our language as follows:

```
Select ?Y, ?Z, ?W
(?X, self :: Biotechnology/backward*, ?Y) AND
(?Y, e-edge/self, ?Z) AND (?Z, e-node, ?W)
```

The first triple binds  $?Y$  to every node reachable from Biotechnology following evolution edges backwards. Then, for each of those nodes, including Biotechnology, the second triple gets all the evolution axes of which the bindings of  $?Y$  are subjects whereas the third triple get the objects of the evolution axes. This query returns, (Biotech,  $\text{evolve}^{-1}$ , Conventional Biotech), (Conventional Biotech,  $\text{join}^{-1}$ , Hybridization), (Conventional Biotech,  $\text{join}^{-1}$ , Fermentation), and (Conventional Biotech,  $\text{join}^{-1}$ , Selective Breeding).

**[Example Query 4]:** Which scientific and engineering concepts and disciplines are related to the emergence of cell cloning? We interpret “related” in our model as being immediate predecessors/successors and “siblings” in the evolution process. That is, from a concept we find its immediate predecessors by following all evolution edges backwards one step. From the results we follow all evolution edges forward on step and we get the original concept and some of its “siblings”. We do the same in the opposite direction following first evolution edges forward one step and then backwards. Based on this notion, we can express the query as follows:

```
Select ?Y, ?Z, ?W
(?X, self :: Cell Cloning, ?Y) AND ?(Y, backward | backward/forward, ?Z)
AND (?Y, forward | forward/backward, ?W)
```

The first triple will just bind Cell Cloning to  $?Y$ . The second triple follows the **detach** edge back to Cloning, and then the **detach** edge forward to Molecular Cloning.

The third triple starts again from Cell Cloning and follows the **join** edge forward to Bioprocess Engineering and then the *detach* edge backwards to Conventional Biotech. All these concepts will be returned by the query.

## 7 Related Work

**Managing time in databases:** Temporal data management [17] has been extensively studied in the relational paradigm [4] (A unified formal framework for studying temporal databases in general and temporal query languages in particular can be found in [18].) For semi-structured data, one of the first models for managing historical information as an extension of the Object Exchange Model (OEM) was introduced in [6]. Versioning schemes for XML were first proposed in [19] and [20], the latter specifically for web data warehouses. These approaches store only the information of the entire document at some point in time and then use edit scripts and change logs to reconstruct versions of the entire document. This may lead to large overheads when processing queries that span multiple versions. In contrast, [21] and [7] maintain a single temporal document from which versions of any document fragment (even single elements) can be extracted directly when needed.  $\tau$ XQuery, a temporal extension to XQuery, was introduced in [22]. In this proposal, temporal queries are translated into XQuery and evaluated by a standard XQuery engine. A survey on temporal extensions to the Entity-Relationship (ER) model is presented in [5].

**Change management in ontologies:** There is a fundamental distinction between an actual update and a revision in knowledge bases [23]. An update brings the knowledge base up to date when the world it models has changed. In contrast, a revision incorporates new knowledge from a world that has not changed. Our evolution framework models updates since it describes how real-world concepts have changed over time. The survey in [24] provides a thorough classification of the types of changes that occur in ontologies. However, their definition of evolution is completely different from ours: they view evolution as a special case of versioning. Moreover, there is no entry in their taxonomy that corresponds to the kind of concept evolution we developed in this work. An approach to model revision in RDF ontologies has been presented in [25]. The first temporal extension to RDF was presented in [8]. Their model is an RDF graph with temporal intervals associated only to triples. Orthogonal to the concept evolution problem is that of terminology evolution which studies how the terms describing the same notion in the domain of discourse have change over time [11]. Similarly to versioning in databases, ontology versioning study the problem of maintaining changes in ontologies by creating and managing different variants of it [26]. The main limitation of such approaches is the same as in databases: large overheads when processing queries that span multiple versions. Closer to our work is the proposal in [27] for modeling changes in geographical information systems (GIS). They use the notion of a change bridge to model how the area of geographical entities (countries, provinces, etc.) evolve. A change bridge is associated with a change point and indicates what concepts become obsolete, what new concepts are created, and how the new concepts overlap with older ones. Since they focus on the GIS domain, they are not able to model causality and different types of evolution involving abstract concepts beyond geographical entities.

## 8 Conclusion

In this work we studied the novel problem of concept evolution, i.e., how the semantics of an entity changes over time. In contrast to temporal models and schema evolution, concept evolution deals with mereological and causal relationships between concepts. Recording concept evolution also allows users to pose queries on the history of a concept. We presented a framework for modeling evolution as an extension of temporal RDF with mereology and causal properties expressed with a set of evolution terms. Furthermore, we presented an extension of nSPARQL that allows navigation over the history of the concepts. Finally, we applied our framework in two real world scenarios, the history of Germany and the evolution of biotechnology, and we showed how queries of interest can be answered using our proposed language.

## References

1. Blakeley, J., Larson, P.A., Tompa, F.W.: Efficiently Updating Materialized Views. In: SIGMOD. (1986) 61–71
2. Lerner, B.S.: A Model for Compound Type Changes Encountered in Schema Evolution. *TODS* **25**(1) (March 2000) 83–127
3. Velegarakis, Y., Miller, R.J., Popa, L.: Preserving mapping consistency under schema changes. *VLDB J.* **13**(3) (2004) 274–293
4. Soo, M.D.: Bibliography on Temporal Databases. *SIGMOD Record* **20**(1) (1991) 14–23
5. Gregersen, H., Jensen, C.S.: Temporal entity-relationship models - a survey. *IEEE Trans. Knowl. Data Eng.* **11**(3) (1999) 464–497
6. Chawathe, S., Abiteboul, S., Widom, J.: Managing historical semistructured data. In: *Theory and Practice of Object Systems*, Vol 5(3). (1999) 143–162
7. Rizzolo, F., Vaisman, A.A.: Temporal XML: modeling, indexing, and query processing. *VLDB Journal* **17**(5) (2008) 1179–1212
8. Gutiérrez, C., Hurtado, C.A., Vaisman, A.A.: Temporal RDF. In: *ESWC*. (2005) 93–107
9. Dong, X., Halevy, A.Y., Madhavan, J.: Reference Reconciliation in Complex Information Spaces. In: *SIGMOD*. (2005) 85–96
10. Palpanas, T., Chaudhry, J., Andritsos, P., Velegarakis, Y.: Entity Data Management in OKKAM. In: *SWAP DEXA Workshop*. (2008) 729–733
11. Tahmasebi, N., Iofciu, T., Risse, T., Niederee, C., Siberski, W.: Terminology evolution in web archiving: Open issues. In: *International Web Archiving Workshop*. (2008)
12. Lenzerini, M.: Data Integration: A Theoretical Perspective. In: *PODS*. (2002) 233–246
13. W3C: RDF vocabulary description language 1.0: RDF Schema. <http://www.w3.org/TR/rdf-schema/> (2004)
14. Dyreson, C.E., Evans, W.S., Lin, H., Snodgrass, R.T.: Efficiently supported temporal granularities. *IEEE Trans. Knowl. Data Eng.* **12**(4) (2000) 568–587
15. Keet, C.M., Artale, A.: Representing and reasoning over a taxonomy of part-whole relations. *Applied Ontology* **3**(1-2) (2008) 91–110
16. Pérez, J., Arenas, M., Gutierrez, C.: nSPARQL: A navigational language for RDF. In: *International Semantic Web Conference*. (2008) 66–81
17. Snodgrass, R.T., Ahn, I.: Temporal Databases. *IEEE Computer* **19**(9) (1986) 35–42

18. Chomicki, J.: Temporal query languages: a survey. In: Proceedings of the 1st International Conference on Temporal Logic, LNAI 827. (1994) 506–534
19. Chien, S., Tsotras, V., Zaniolo, C.: Efficient management of multiversion documents by object referencing. In: Proceedings of the 27th International Conference on Very Large Data Bases, Rome, Italy (2001) 291–300
20. Marian, A., Abiteboul, S., Cobena, G., Mignet, L.: Change-centric management of versions in an XML warehouse. In: Proceedings of the 27th VLDB Conference, Rome, Italy (2001) 581–590
21. Buneman, P., Khanna, S., Tajima, K., Tan, W.: Archiving scientific data. In: Proceedings of the 2002 ACM SIGMOD International Conference on Management of Data, Madison, USA (2002) 1–12
22. Gao, C., Snodgrass, R.: Temporal slicing in the evaluation of XML queries. In: Proceedings of the 29th International Conference on Very Large Data Bases, Berlin, Germany (2003) 632–643
23. Katsuno, H., Mendelzon, A.O.: On the difference between updating a knowledge base and revising it. In: International Conference on Principles of Knowledge Representation and Reasoning. (1991) 387–394
24. Flouris, G., Manakanatas, D., Kondylakis, H., Plexousakis, D., Antoniou, G.: Ontology change: classification and survey. *Knowledge Eng. Review* **23**(2) (2008) 117–152
25. Konstantinidis, G., Flouris, G., Antoniou, G., Christophides, V.: On RDF/S ontology evolution. In: SWDB-ODDIS. (2007) 21–42
26. Klein, M.C.A., Fensel, D.: Ontology versioning on the semantic web. In: SWWS. (2001) 75–91
27. Kauppinen, T., Hyvönen, E.: Modeling and reasoning about changes in ontology time series. In: *Ontologies: A Handbook of Principles, Concepts and Applications in Information Systems*. (2007) 319–338