

Chapter 1

Dealing with Structure Heterogeneity in Semantic Collaborative Information Systems

Eva Zangerle

University of Innsbruck, Austria

Wolfgang Gassler

University of Innsbruck, Austria

ABSTRACT

The creation of content within semistructured, collaborative information systems imposes the problem of having to deal with very heterogeneous schemata. This is due to the fact that the semistructured paradigm does not restrict the user in his choice of nomenclatures for the data he intends to store within the information system. As many users participate in the creation of data, the structure of this data is very heterogeneous. In this chapter the authors discuss two main movements that aim at dealing with heterogeneity. The first approach is concerned with efficiently avoiding structure heterogeneity within collaborative information systems by providing the users with suitable recommendations for an aligned schema during the insertion process. The second approach is mainly focussing on overcoming structure heterogeneity by providing efficient means for querying heterogeneous data.

INTRODUCTION

Most online, collaborative information systems, such as wiki systems, provide means to easily add, modify and delete information, which does not have to adhere to any predefined schema or

structure. In contrast, traditional (relational) databases are strictly-structured and enforce the user to store information in a predefined schema. Such structured data stores provide the big advantage of structured access, which enables complex query capabilities. Traditional wiki systems only support full-text search which is not feasible for complex

DOI: 10.4018/978-1-4666-0894-8.ch001

queries such as “Which Austrian cities have more than 10.000 inhabitants and have a female mayor who has a doctoral degree?” Nevertheless, wiki systems are able to cope with very large amount of collaboratively created information with very heterogeneous structures and schemata.

Weikum et al. (2009) observed that modern information systems have to be able to support both structured and unstructured data to combine the advantages of both worlds and be able to answer such complex questions. This need of combination initiates the development of collaborative, semistructured information systems. They provide mechanisms for the combination of both unstructured and structured storage of data. Semistructured data features a structure without having to specify a fixed schema. As this paradigm does not restrict the user and the used schema at all, the massive collaborative creation and editing of content by hundreds or thousands of users obviously leads to the usage of very heterogeneous schemata and structures in collaborative environments. Even Wikipedia, which has a very committed community dealing with heterogeneity, is also not able to avoid heterogeneity within its schema.

In the following sections we discuss the problem of heterogeneity in semistructured information systems and show approaches which are able to deal with heterogeneous schemata, data and the collaborative paradigm of creating and managing knowledge and information.

Schema and Heterogeneity

Modern collaborative Information Systems mostly use the semistructured paradigm, as it features the possibility to structure information without having to adhere to a predefined schema. The most popular example of a semistructured data format is RDF which is often used in the underlying storage layer of semantic, collaborative information systems. RDF consists of triples with the form $\langle \text{subject}, \text{predicate}, \text{object} \rangle$. For example

information about the subject Albert Einstein can be stored by using the triples shown in Listing 1.

Listing 1. Semistructured description of Albert Einstein

```
<AlbertEinstein,name,Albert Einstein>
<AlbertEinstein,born,1897/03/14>
<AlbertEinstein,bornIn,Ulm>
<AlbertEinstein,wonAward,Nobel Prize>
<AlbertEinstein,wonAward,Max Planck Medal>
```

RDF distinguishes between URIs which describe resources and literals for specifying values. For reasons of simplification, the URI prefixes are omitted in the example in Listing 1. RDF can also be represented as a graph $G = (N, L, E)$ where N contains all nodes. The edges of the graph are defined in the set E with $e(n_1, n_2)$ and $n_1, n_2 \in N$. The possible labels of edges are denoted by L . In the graph representation, subjects and objects are modelled as nodes. The predicates are modelled as labels and the triple $\langle s_i, p_i, o_i \rangle$ itself is defined by an edge $e(s_i, o_i)$ with the label p_i .

The schema or structure of a graph is defined similar to a classical relational database schema. A database schema is defined by the columns - also called attributes - of a relation. Each row also called record of a relation has to conform to a predefined schema. As semistructured data or RDF data are not restricted in any way, each record can consist of arbitrary many attributes. In the context of RDF a record is called subject and attributes are called predicates. The schema or structure of a RDF subject S_i is defined by the used attributes

$$Schema_{S_i} = \{p \mid p \in L, e(S_i, n) \text{ with label } p\}$$

In contrast to classical relational databases each record can constitute its own schema. Exactly this feature of RDF is one of the most important

18 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/dealing-structure-heterogeneity-semantic-collaborative/65684

Related Content

Search Engine-Based Web Information Extraction

Gijs Geleijnse (2009). *Semantic Web Engineering in the Knowledge Society* (pp. 208-241).

www.irma-international.org/chapter/search-engine-based-web-information/28854

Social Presence and User-Generated Content of Social Media in China

Rui Sun and Hong Xue-Jiao (2019). *International Journal on Semantic Web and Information Systems* (pp. 35-47).

www.irma-international.org/article/social-presence-and-user-generated-content-of-social-media-in-china/227351

Developing a Web-Based Cooperative Environment to Software Project Development

Seyed Morteza Babamir (2012). *Collaboration and the Semantic Web: Social Networks, Knowledge Networks, and Knowledge Resources* (pp. 246-270).

www.irma-international.org/chapter/developing-web-based-cooperative-environment/65696

Integration of Semantics Into Sensor Data for the IoT: A Systematic Literature Review

Besmir Sejdiu, Florije Ismaili and Lule Ahmed (2020). *International Journal on Semantic Web and Information Systems* (pp. 1-25).

www.irma-international.org/article/integration-of-semantics-into-sensor-data-for-the-iot/264161

Family History Information Exchange Services Using HL7 Clinical Genomics Standard Specifications

Amnon Shabo and Kevin S. Hughes (2005). *International Journal on Semantic Web and Information Systems* (pp. 44-67).

www.irma-international.org/article/family-history-information-exchange-services/2814