

# BauDenkMalNetz – Creating a Semantically Annotated Web Resource of Historical Buildings

Anca Dumitrache and Christoph Lange

Computer Science, Jacobs University Bremen, Germany {a.dumitrache,ch.lange}@jacobs-university.de

**Abstract.** BauDenkMalNetz (“listed buildings web”) deals with creating a semantically annotated website of urban historical landmarks. The annotations highlight the most relevant information about the landmarks (e.g. the buildings’ architects, architectural style or construction details), for the purpose of extended accessibility and smart querying. BauDenkMalNetz is based on a series of touristic books on architectural landscape. After a thorough analysis on the requirements that our website should provide, we processed these books using automated tools for text mining, which led to an ontology that represents the vocabulary necessary to express all relevant architectural and historical information. In preparation of publishing the books on a website powered by this ontology, we analyze how well Semantic MediaWiki and the RDF-aware Drupal 7 content management system satisfy our requirements for browsing the content in a semantic way.

## 1 Motivation

The architectural landscape of a city is generally made up not just of well-established landmarks, but of historical buildings with a rich cultural background that lie outside the mainstream touristic circuit. People wanting to explore the more personal and less-known places of a city have little access to information about these hidden architectural gems and the stories behind them, even though all required data on historical buildings in Germany has been meticulously collected by the offices for historical monuments (Denkmalämter). However, this data has generally not been published in an easily accessible way. Existing databases and form-based search facilities are often tedious to browse through.<sup>11</sup>

In Bremen, an effort to collect this information and present it to the general public was made by the publisher Nils Aschenbeck, who released a series of guide books [AW09] about the city. However, for the moment, these books are only accessible in printed format. By making use of these books, BauDenkMalNetz (German for “listed buildings web”) proposes a way of discovering Bremen’s architectural landscape that is suited for the tech-savvy tourist.

EdNote(1)

<sup>11</sup> See, for example, <http://194.95.254.61/denkmalpflege/index.htm>.

<sup>1</sup> EDNOTE: CL: maybe turn into a bib entry

## 2 Transitioning from Written Text to Digital Media

General-interest media publications usually make use of a concrete set of concepts, that relate to one particular subject area, and thus can be reduced to a strict vocabulary. By grouping media publications into conceptual domains, the semantic structures can be reused across publications, so that they can be linked across the Semantic Web, creating a network of interconnected publications and information. Our purpose with developing the BauDenkMalNetz data representation was creating a general model for our particular conceptual domain. Therefore, we have conceptualized our knowledge bases in such a way that the conceptualization can be reused for representing other texts on the topic of historical landmarks.

### 2.1 Building an Ontology

The publications that lie at the basis of our work with BauDenkMalNetz are stored in simple HTML files. There is a file for each individual building, with pictures associated to each file, and information like the name of the architect being highlighted. Four books have been published thus far [AW09], with more than one hundred buildings being described in total, which makes for a large amount of data to be processed. In order to enable enhanced browsing and querying, the data on Bremen's historical buildings needs to be organized, and the proper semantic metadata needs to be put in place. For this purpose, we have developed the BauDenkMalNetz ontology, a formal representation of the metadata vocabulary on historical buildings and related concepts, together with the relations among them. The ontology has been formalized and implemented in OWL, and was engineered in the stages specified by the METHONTOLOGY [FLGPJ97] ontology engineering methodology.

**Scenario** An example scenario of interacting with a publication backed by the BauDenkMalNetz ontology involves a tourist, who wishes to work out an itinerary for visiting the city of Bremen. For this purpose, she needs to be able to browse through a particular neighborhood, by filtering the buildings based on their addresses. Furthermore, she is interested only in visiting those buildings that were built in the nineteenth century. Then she finds one particular architect that she is familiar with, and she wants to add all of his buildings to her itinerary. Finally, during her visit, she will want to stop at each individual building and read up on its history, like the years between it was built, and what famous people had been living there.

**Requirements** Based on this scenario, we have identified a list of requirements that the BauDenkMalNetz ontology needs to meet in order for the data to be easily accessible:

- *buildings* need to be represented as uniquely identified entities, which will be mapped to individual pages of the website; any knowledge represented using the BauDenkMalNetz ontology needs to be interconnected, with the building entity as the central point of the representation;
- information on the *physical address* and *neighborhood* needs to be available for every building;
- the *architect* and the *architectural style* of a building have to be highlighted when that information is available;
- the *time* and *timespan* over which a building was built has to be specified for individual entries.

A more general requirement that the BauDenkMalNetz website needs to address is browsing from one building to another. This could be done via information on the buildings' physical location (e.g. they are on the same street), or based on characteristics that they share (e.g. they were built by the same person).

**Text Analysis** Starting from these requirements and based on the original touristic guides on which the website will be based, we identified the vocabulary that relates to historical buildings, by employing **n-gram models** [MS99] to find the most likely occurrences of word groupings. The results of this analysis were used in the conceptualization phase of the BauDenkMalNetz ontology.

The first step that enabled us to process the text was removing the unnecessary HTML metadata, and stripping it down to a plain-text format. The text is written in German; we needed to normalize it to plain ASCII characters, as the German-specific special characters seemed to interfere with the script used to analyze it. We made use of the **LaMaPUn** [GJA+09] Perl library for processing the text. We used a list of the most frequent German stop words in order to filter out the information that was not meaningful for the domain vocabulary.

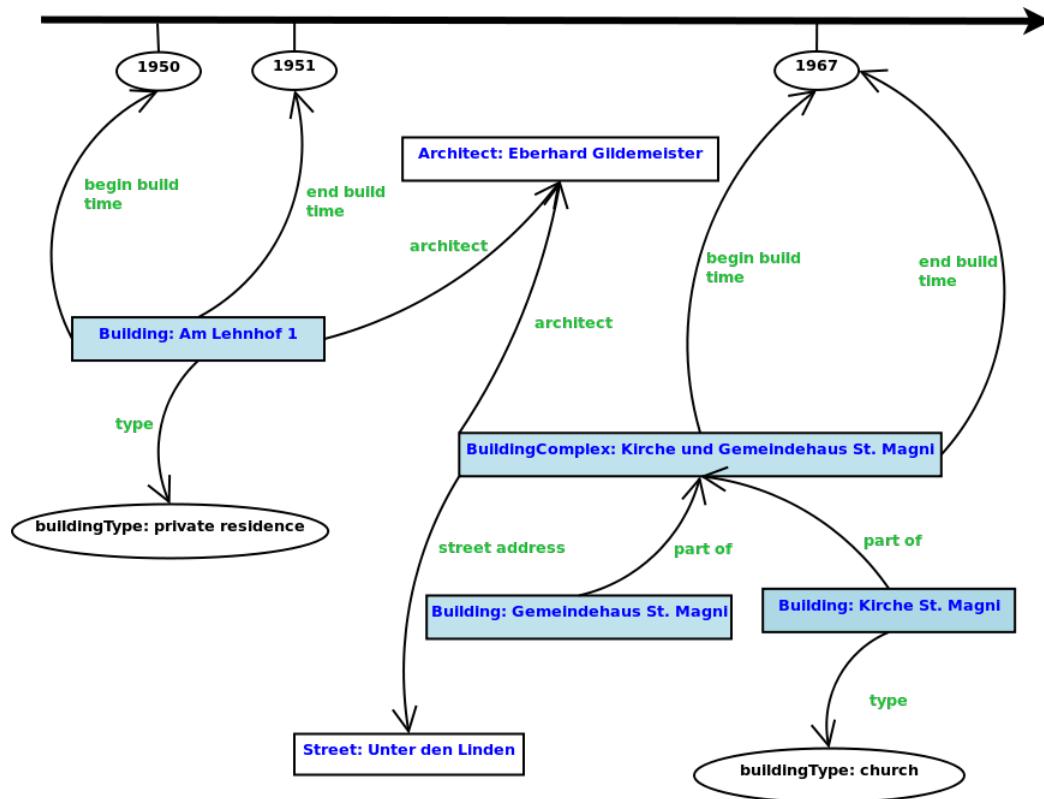
An n-gram model refers to a probabilistic model that, given the first  $n - 1$  words in a sentence, will predict the  $n^{\text{th}}$  word. By creating such a model, we were able to get a first impression of what concepts will be mapped to resources and properties in our RDF model of the metadata.

We analyzed series of 1 to 4-gram models. The script recognized over 600 possible groupings of words that are likely to occur together. Over 500 of these groups had a likelihood coefficient larger than 2. This coefficient is computed by having the number of incidences of the words in the group together divided by the sum of individual incidences outside of the group.

The text analysis made apparent some clear trends. Most of the likely groups of words that appeared together referred to one of the following categories: *physical buildings* (e.g. Bahnhof (*train station*) Sankt Magnus, Kirche (*church*) Sankt Magni), *personal names* (e.g. Rudolf Alexander Schroeder), *physical addresses* (e.g. Leuchtenburger Strasse (a *street*), Am Bahnhof Sankt Magnus) and *building features* (e.g. Bungalow, Turm (*tower*)). By identifying these categories, we got a first impression of what resources we need to define for our ontology.

**Conceptualization** Based on this analysis, and according to the requirements identified in the previous section, we conceptualized the following entities to be represented in the BauDenkMalNetz ontology (concepts underlined, relations in *italics*):

- building – a resource identifying a particular building;
- building part – a subconcept of the building entity (e.g. tower, annex);
- building complex – a composite consisting of several building entities;
- building type – different types of constructions (e.g. church, hospital);
- *address* – the physical location of a building;
- *architect* – the person or group of people that have designed the building;
- *inhabitant* – famous person that has lived in that building;
- *year* – *when a building was built*; can refer to the year when *construction began, ended, or both*.



**Fig. 1.** A fragment of the BauDenkMalNetz ontology

**Alignment to Other Ontologies** The linked-data community [Hea+] advocates the reuse of knowledge models and vocabularies, in order to achieve interoperability across the Web. Indeed, there already exist various ontologies that model some of the relevant knowledge about historical buildings, out of which we found the following ones relevant for aligning with the BauDenkMalNetz ontology:

- The **GeoNames** [Geo] ontology models geospatial semantic information. In particular, it assigns to individual locations on the globe a unique URI with a corresponding RDF web service. For our purposes, it can be used to uniquely identify each historical building based on its coordinates. Reusing this ontology brings the added advantage of explicitly specifying the geolocation of a building, which allows for easier integration with web mapping services.
- The **CIDOC CRM** [Cid] ontology represents the detailed scientific documentation of cultural heritage objects, which include historical monuments. By aligning our ontology to CIDOC CRM, we can formulate a full description of the historical information related to a building (e.g. the architectural style of the monument, the official sources which document the monument etc.).

## 2.2 Publishing in a Semantic Content Management System

For deploying BauDenkMalNetz as a semantic publication, we have so far established requirements and analyzed how well two semantic content management systems satisfy these requirements: **Semantic MediaWiki** (SMW [Sem]) and **Drupal 7** [Dru].

**Requirements** Based on the scenario discussed in the previous section, we have also analyzed the requirements that our website needs to provide. Digitally representing publications means that the BauDenkMalNetz web portal needs to build on the use cases of the written text that lies at its core, and enhance them with semantic browsing and querying capabilities that will provide for a better user experience. Therefore, a suitable content management system for deploying BauDenkMalNetz should offer the following functionality:

1. the possibility of integrating RDF triples, and at least a minimum of ontology support;
2. support for querying the RDF content of the website (e.g. by using SPARQL);
3. browsing based on the semantic metadata;
4. extensible publishing support for:
  - (a) people, through enabling PDF and HTML exporting;
  - (b) machines, by interlinking the publications across the Web, according to linked data principles;
5. the possibility of importing large amounts of text into the system.

**Semantic MediaWiki** SMW [Sem] was built as an extension of MediaWiki, the well known wiki engine which powers Wikipedia. It provides enhanced features for browsing and organizing its contents as a result of the added semantic annotations to its text. We built the first BauDenkMalNetz prototype using SMW [DLK+10].

Our motivation for using SMW in deploying the initial version of our web portal was its suitability for rapidly creating a working prototype (cf. [BDH+09]). SMW allows easily adding and editing of the necessary data and metadata available on historical buildings, in keeping with requirements 1 and 3. New information could be easily incorporated and linked to the already existing data by making use of SMW page creation and editing tools. At the same time, the structure of the metadata (i.e. the ontology representing the data) could be easily modified, simply by adding in-text annotations.

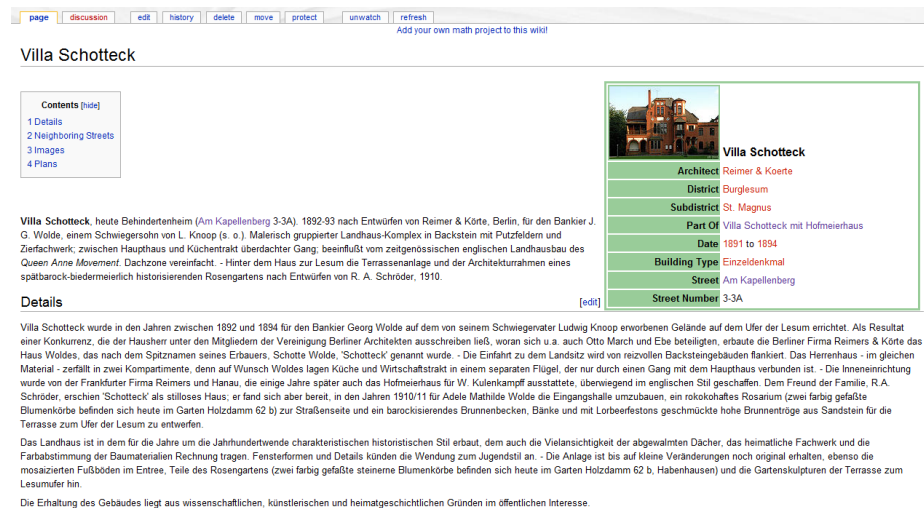


Fig. 2. Screenshot of the SMW prototype

However, when further assessing requirement 1, we found that the conceptual model of our metadata was less obvious and never explicitly formalized, as the ontology, to which the texts adhere, is not necessarily specified explicitly in SMW, but rather implied from the annotations done directly on the text. In this case, alignment to other similar ontologies (in keeping with the linked-data philosophy of reuse) is still possible, yet it is rendered more difficult by the lack of an explicit formal definition of the ontology.

Requirement 5 was also not addressed by our prototype. SMW provides some tools suited for database import, however the texts we want to analyze are stored in simple HTML files. The volume of data that needs to be processed makes it almost impossible to have the texts annotated manually, like we did for

building the prototype, while also making BauDenkMalNetz rather suited for the employment of natural language processing techniques in order to get the needed semantical annotations.

**Drupal 7** As our goal is to publish existing content, rather than creating new content in a collaborative way, we also considered **Drupal** [Dru], a rather traditional content management system. Given the BauDenkMalNetz documents collection and our ontology, we have so far analyzed Drupal's features w.r.t. the requirements established above. Deploying BauDenkMalNetz in Drupal remains to be done in spring 2011.

The latest version 7 of Drupal provides an RDF API [CDC+09], which will enable us to easily integrate our OWL ontology into the website, by using Drupal's taxonomy XML import feature. Afterwards, the keywords pertaining to each resource will be added to the system, and mapped to the corresponding classes and properties in the ontology. For printed media, where a particular text usually does not undergo much change after being published, the advantage that Drupal brings is that, as the structure of the text is already known, its conceptualization can be set as the core of the website via the RDF API even before the website is deployed.

**Comparison** When comparing SMW to Drupal, we have encountered some drawbacks of SMW that led us to reconsider our approach. The flexibility and agility of SMW were not of a particular advantage in our setting. The publication sources are imported from external sources, and therefore we are not interested in MediaWiki's collaboration support. The ontology and its connections to other ontologies are, for now, created just by us, but they are not evolved or extended dynamically by a community – therefore we are not interested in giving write access to the ontology via the content management system. We rather prefer having a clear conceptual model of the metadata from the beginning. Drupal supports the initial import of such an ontology before importing the content and thus is suited for managing annotations to publications that have already existed before.

### 3 Development and Evaluation Plan

During spring 2011, we will continue developing the BauDenkMalNetz website in Drupal, starting with importing the ontology. Next, we will import the texts of the tourist guides, with the keywords in the vocabulary highlighted in the resource's pages. We will make a semantic search available, based on these key concepts, achieved through Drupal's taxonomy feature. Also, for increased functionality, we will add a geospatial aspect to the semantic navigation by utilizing the Google Maps API [Goo]. For even more advanced querying features, we are considering to make use of the **XSPARQL** [AKK+08] query language. XSPARQL combines the XML query language XQuery with the RDF query

language SPARQL, which allows for obtaining XML results for queries over the semantic metadata of our website and, in future, interlinked websites. By selecting from a list of available queries, tourists will be able to create personalized guides of historical buildings.

For evaluating the usability of the BauDenkMalNetz website, existing methods for evaluating (semantic) digital libraries [FTA+07; Kru09] are applicable. A group of test-users will navigate through the website, providing feedback based on *usability* (of the content management system with our extensions) and *usefulness* (of the content, in the way our system publishes it). The users will provide feedback on how easy/difficult it is to find a particular building, by querying the system based on a criteria of their own choosing (e.g. location, architectural style etc.), and also about how they managed to find their way from one particular building to another, based on a common characteristic. They will also be asked to provide their input on how accurate the query results are in relation to what they were expecting to find, and also about the informative character of individual buildings' pages. Based on this evaluation, we will assess the user-friendliness of the website and consider possible improvements. A first release of BauDenkMalNetz, adapted according to the results of an initial evaluation round, is expected in early May.

## 4 Related Work

MANTIC [MPV10] is a project similar to BauDenkMalNetz, that represents data on cultural heritage sites of the city of Milan, that was gathered from historical sources and publications. It has the CIDOC CRM ontology at its core, which it uses in order to store information about the archeology of the city. This information is then incorporated into the Google Maps API, making for an easy to use application for browsing Milan's historical landmarks, that is quite similar in scope to our work.

Unlike BauDenkMalNetz, MANTIC deals with historical sources, which comprise a great variety of publications, written in different styles and over a long period of time. MANTIC provides a good example of how CIDOC CRM can be reused for representing historical landmarks, however, since the sources MANTIC deals with are so disjointed, identifying a common vocabulary for them is more difficult, and therefore no special ontology that deals primarily with historical buildings was devised.

## 5 Conclusion and Further Work

After assessing in which ways traditional printed publications on historical landmarks can be enhanced by transposing them in a digital format and enriched with semantic annotations, we devised the BauDenkMalNetz ontology, by analyzing its requirements and processing the texts that were made available to us by using text mining techniques. In keeping with linked data principles, we aligned our ontology to other existing representations that relate to our specific



domain, like CIDOC CRM and GeoNames. Once we determined the structure of our metadata, we compared how different content management systems (SMW and Drupal 7) satisfy the requirements for deploying the BauDenkMalNetz website. As Drupal provides a more rigorous way of declaring a conceptual model, which is more suitable for digital publications, we have chosen it as the medium in which our web portal will be developed.

Once finished, the BauDenkMalNetz website will provide a comprehensive and easy-to-use guide to the city of Bremen, and possibly even help boost the touristic appeal of Bremen. A possible enhancement to the resource will be creating a mobile version of the website, so that tourists can create virtual itineraries that they can access on the go. However, the scope of our work is not limited to Bremen. Both the ontology and the vocabulary are general enough to adapt in order to represent any touristic publication guide on historical landmarks.

## Acknowledgments

The authors would like to thank Deyan Ginev for help with the LaMaPUn library, and Lin Clark for help with assessing Drupal 7.

## References

- [AKK+08] W. Akhtar, J. Kopecký, T. Krennwallner, and A. Polleres. “XS-PARQL: Traveling between the XML and RDF worlds – and avoiding the XSLT pilgrimage”. In: *The Semantic Web: Research and Applications*. 5<sup>th</sup> European Semantic Web Conference (ESWC) (Tenerife, Spain, June 1–5, 2008). Ed. by S. Bechhofer, M. Hauswirth, J. Hoffmann, and M. Koubarakis. Lecture Notes in Computer Science 5021. Springer Verlag, 2008.
- [AW09] N. Aschenbeck and I. Windhoff. *Landhäuser und Villen in Bremen*. Bremen: Aschenbeck Verlag, 2009.
- [BDH+09] J. Bao, L. Ding, R. Huang, P. Smart, D. Braines, and G. Jones. “A Semantic Wiki based Light-Weight Web Application Model”. In: *Proceedings of the 4<sup>th</sup> Asian Semantic Web Conference*. 2009, pp. 168–183.
- [CDC+09] S. Corlosquet, R. Delbru, T. Clark, A. Polleres, and S. Decker. “Produce and Consume Linked Data with Drupal!” In: *The Semantic Web – ISWC 2009*. 8<sup>th</sup> International Semantic Web Conference (ISWC). Ed. by A. Bernstein, D. R. Karger, T. Heath, L. Feigenbaum, D. Maynard, E. Motta, and K. Thirunarayan. LNCS 5823. Springer Verlag, Oct. 2009.
- [Cid] *The CIDOC Conceptual Reference Model*. URL: <http://cidoc.ics.forth.gr> (visited on 2010-03-07).

- [DLK+10] A. Dumitrache, C. Lange, M. Kohlhase, and N. Aschenbeck. “Prototyping a Browser for a Listed Buildings Database with Semantic MediaWiki”. In: *5<sup>th</sup> Workshop on Semantic Wikis*. (Hersonissos, Crete, Greece, May 31, 2010). Ed. by C. Lange, J. Reutelshöfer, S. Schaffert, and H. Skaf-Molli. CEUR Workshop Proceedings 632. Aachen, 2010. URL: <http://kwarc.info/clange/pubs/semwiki2010-baudenkmalnetz.pdf>.
- [Dru] *Drupal.org – Community plumbing*. web page at <http://drupal.org>. URL: <http://drupal.org>.
- [FLGPJ97] M. Fernández-López, A. Gómez-Pérez, and N. Juristo. “METHONTOLOGY: from Ontological Art towards Ontological Engineering”. In: *Proceedings of the Fourteenth National Conference on Artificial Intelligence AAAI-97*. (Stanford, USA, Mar. 1997). MIT Press, 1997, pp. 33–40.
- [FTA+07] N. Fuhr, G. Tsakonas, T. Aalberg, M. Agosti, P. Hansen, S. Kapidakis, C.-P. Klas, L. Kovács, M. Landoni, A. Micsik, C. Papatheodorou, C. Peters, and I. Sølvberg. “Evaluation of digital libraries”. In: *International Journal of Digital Libraries* 8 (2007), pp. 21–38.
- [Geo] *GeoNames*. URL: <http://www.geonames.org> (visited on 2010-04-23).
- [GJA+09] D. Ginev, C. Jucovschi, S. Anca, M. Grigore, C. David, and M. Kohlhase. “An Architecture for Linguistic and Semantic Analysis on the arXMLiv Corpus”. In: *Applications of Semantic Technologies (AST) Workshop at Informatik 2009*. 2009. URL: [http://www.kwarc.info/projects/lamapun/pubs/AST09\\_LaMaPUn+appendix.pdf](http://www.kwarc.info/projects/lamapun/pubs/AST09_LaMaPUn+appendix.pdf).
- [Goo] *Google Maps*. URL: <http://maps.google.com> (visited on 2011-01-10).
- [Hea+] T. Heath et al. *Linked Data – Connect Distributed Data across the Web*. URL: <http://linkeddata.org> (visited on 2010-06-11).
- [Kru09] S. R. Kruk. “Semantic Digital Libraries. Improving Usability of Information Discovery with Semantic and Social Services”. PhD thesis. National University of Ireland, Galway, 2009.
- [MPV10] G. Mantegari, M. Palmonari, and G. Vizzari. “Rapid Prototyping a Semantic Web Application for Cultural Heritage: The Case of MANTIC”. In: *The Semantic Web: Research and Applications (Part II)*. 7<sup>th</sup> Extended Semantic Web Conference (ESWC) (Hersonissos, Crete, Greece, May 30–June 3, 2010). Ed. by L. Aroyo, G. Antoniou, E. Hyvönen, A. ten Teije, H. Stuckenschmidt, L. Cabral, and T. Tudorache. Lecture Notes in Computer Science 6089. Springer Verlag, 2010.
- [MS99] C. D. Manning and H. Schütze. “Statistical Inference: n-gram Models over Sparse Data”. In: *Foundations of Statistical Natural Lan-*

[Sem] *guage Processing*. Ed. by Name. Cambridge, Massachusetts: MIT Press, 1999. Chap. 6.  
*Semantic MediaWiki*. URL: <http://semantic-mediawiki.org> (visited on 2010-03-04).