

MultimediaN E-Culture demonstrator

Guus Schreiber¹, Alia Amin², Mark van Assem¹, Victor de Boer³, Lynda Hardman^{2,5}, Michiel Hildebrand², Laura Hollink¹, Zhisheng Huang², Janneke van Kersen⁴, Marco de Niet⁴, Borys Omelayenko¹, Jacco van Ossenbruggen², Ronny Siebes¹, Jos Taekema⁴, Jan Wielemaker³, and Bob Wielinga³

¹ Vrije Universiteit Amsterdam (VU), Computer Science, <http://www.cs.vu.nl>

² Center for Math. and Computer Science CWI, Amsterdam, <http://www.cwi.nl/ins2>

³ Universiteit van Amsterdam (UvA), HCS Lab, <http://hcs.science.uva.nl/>

⁴ Digital Heritage Netherlands (DEN), The Hague, <http://www.den.nl>

⁵ Technical University Eindhoven (TU/e), <http://w3.win.tue.nl/en/>

1 Introduction

The main objective of the MultimediaN E-Culture project is to demonstrate how novel semantic-web and presentation technologies can be deployed to provide better indexing and search support within large virtual collections of cultural-heritage resources. The architecture is fully based on open web standards, in particular XML, SVG, RDF/OWL and SPARQL. One basic hypothesis underlying this work is that the use of explicit background knowledge in the form of ontologies/vocabularies/thesauri is in particular useful in information retrieval in knowledge-rich domains.

This paper gives some details about the internals of the demonstrator. The online version of the demonstrator can be found at:

<http://e-culture.multimedian.nl/demo/search>

Readers are encouraged to first take a look at the demonstrator before reading on. As a teaser we have included a short description of basic search facilities in the next section. We suggest you consult the tutorial (linked from the online demo page) which provides a sample walk-through of the search functionality. Make sure your browser has adequate SVG support.⁶

Please note that this is a product of an ongoing project. Visitors should expect the demonstrator to change⁷. We are incorporating more collections and vocabularies and are also extending the annotation, search and presentation functionality.

⁶ The current version of the demonstrator runs under Firefox version 1.5.0.4 with the Adobe SVG plugin (v 6.0 38363, see the demonstrator FAQ for installation instructions) and has been tested on Windows, Macintosh and Linux. Support for Internet Explorer is planned for future versions of the demo. Firefox 2 is expected to make the plugin installations unnecessary (you can try the beta-release). As a project we are committed to web standards (such as SVG) and are not willing to digress to (and spend time on) special-purpose solutions.

⁷ The project has a duration of 4 years and is at the time of writing 18 months underway.

2 A Peek at the Demonstrator

Figure 1 shows a query for Art Nouveau. This query will retrieve images that are related to Art Nouveau in some way. The results shown in the figure are ‘created by an artist with a matching style’. So, these images are paintings by artists who have painted in the Art-Nouveau style, but the style is not part of the metadata of the image itself. This may retrieve some paintings which are not really Art Nouveau, but it is a reasonable strategy if there are no (or only few) images directly annotated with Art Nouveau. We view the use of such indirect semantic links as a potential for semantic search (for more details on path search in the demonstrator see Section 8).

The lower part of the figure shows a listing of painters who are known to have worked in the Art-Nouveau style. The time line indicates years in which they have created art works (you can click on them to get information).

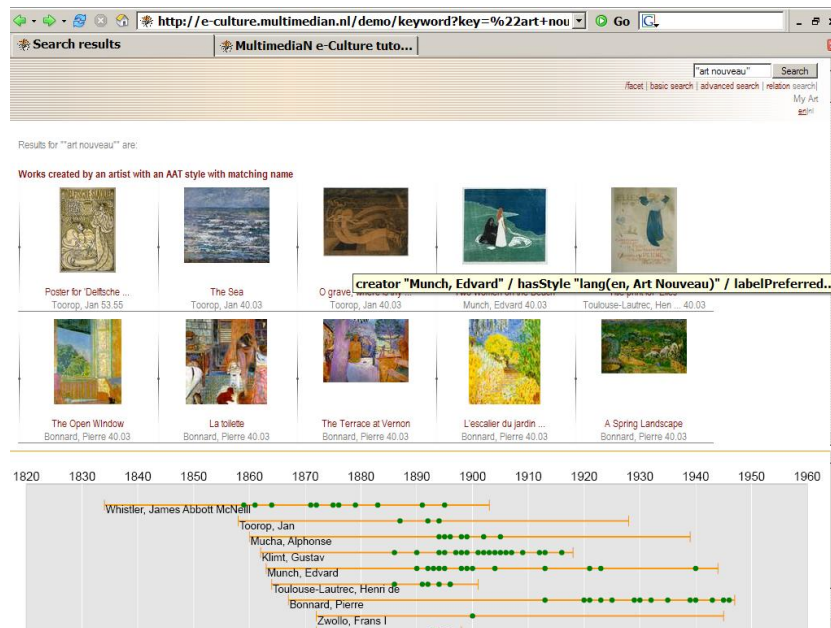


Fig. 1. Results of query for Art Nouveau

Images have annotations in which terms from various thesauri are used. Figure 2 shows the information a user gets when selecting such an indexing term, here **Gold** material from the Art & Architecture Thesaurus. We also show images that have been annotated with this indexing term (or semantically related terms).


These are some basic search- and presentation functions. See the online demo for information about more search options, such as time-based search and faceted

```

Description:
aat:related to    aat:precious metal, aat:goldwork, aat:elements;
vp:descriptiveNote Use for the pure metallic element having symbol Au and atomic number 79; soft, heavy, chemically inactive, yellow metal, considered as a precious metal since ancient times, and
                  serving in many cultures as the basis of material trade values. Use also for this metal as processed and formed, usually in combination with other substances, to make various objects
                  and materials.;
vp:id            300011021;
vp:labelNonPreferred Au;
vp:labelPreferred gold;
vp:preferred parent aat:gold and gold alloy;
rdf:type         aat:Concept;

Value:
aat:precious metal, aat:goldwork, aat:elements;
aat:shell gold;

```



```

aat:related to
vp:preferred parent
vra:Material:Medium

```

Fig. 2. Information about the indexing term **Gold**, showing also images that are related to **Gold**

search. We also have an experimental search function for finding the semantic relations between two URIs, e.g. for posing the question “How are Van Gogh and Gauguin related?”. In fact, this leads to a whole avenue of new search possibilities and related issues with respect to which semantic paths are most relevant, which we hope to explore in more detail in the coming years.

3 Technical Architecture

The foundation of the demo is formed by SWI-Prolog and its (Semantic) Web libraries (for detailed information, see [1, 2]). SPARQL-based access is a recent feature. The *Application Logic* module defines searching and clustering algorithms using Prolog as query language, returning the results as Prolog Herbrand terms. The *Presentation Generation* module generates web documents from the raw answers represented as Herbrand terms.

From the user perspective, the architecture provides (i) annotation facilities for web resources representing images, and (ii) search and presentation/visualization facilities for finding images.

4 Vocabularies

Currently, the demonstrator hosts four thesauri, namely the three Getty vocabularies⁸, i.e., the Art & Architecture Thesaurus (AAT), Union List of Artists Names (ULAN) and the Thesaurus of Geographical Names (TGN), as well as

⁸ http://www.getty.edu/research/conducting_research/vocabularies/

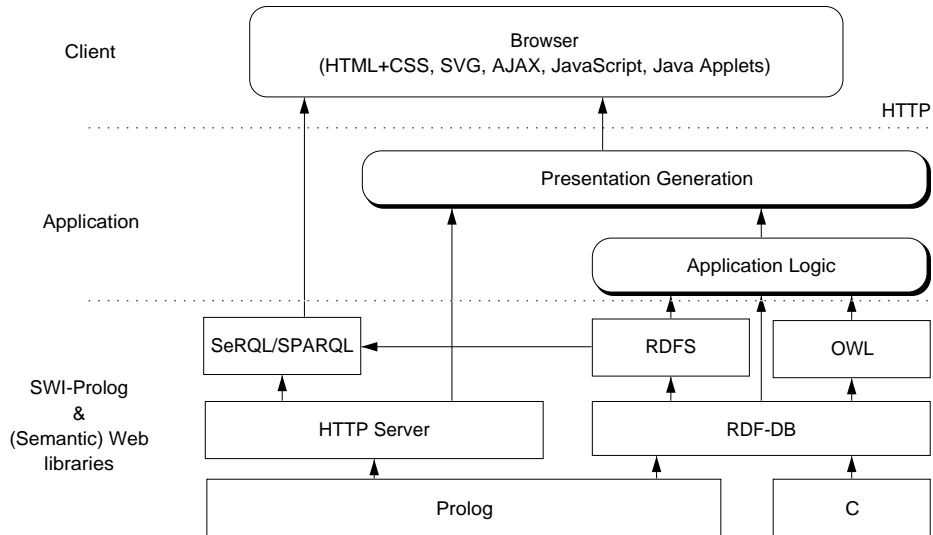


Fig. 3. Technical architecture of the demonstrator

the lexical resource WordNet, version 2.0. The Getty thesauri were converted from their original XML format into an RDF/OWL representation using the conversion methods principles as formulated in [3]. The RDF/OWL version of the data models is available online⁹. The Getty thesauri are licensed¹⁰.

The RDF/OWL conversion of WordNet is documented in a publication of the W3C Semantic Web Best Practices and Deployment Working Group [4]. It is an instructive example of the issues involved in this conversion process, in particular the recipes for publishing RDF vocabularies [5].

The architecture is independent of the particular thesauri being used. We are currently in the process of adding the Dutch version of AAT, amongst others to support a multi-lingual interface. Integration of other (multi-lingual) thesauri is planned.

Using multiple vocabularies is a baseline principle of our approach. It also raises the issue of alignment between the vocabularies. Basically, semantic interoperability will increase when semantic links between vocabularies are added. Within the Getty vocabularies one set of links is systematically maintained: places in ULAN (e.g., place of birth of an artist) refer to terms in TGN. Within the project we are adding additional sets of links. One example is links between art styles in AAT (e.g. “Impressionism”) and artists in ULAN (e.g., “Monet”). De Boer [6] has worked on deriving these semi-automatically from texts on art history.

⁹ <http://e-culture.multimedien.nl/resources/>

¹⁰ The partners in the project have acquired licenses for the thesauri. People using the demonstrator do not have access to the full thesauri sources, but can use them to annotate and/or search the collections.

5 Annotation template

For annotation and search purposes the tool provides the user with a description template derived from the VRA 3.0 Core Categories [7]. The VRA template is defined as a specialization of the Dublin Core set of metadata elements, tailored to the needs of art images. The VRA Core Categories follow the “dumb-down” principle, i.e., a tool can interpret the VRA data elements as Dublin Core data elements.¹¹

6 Collection data and metadata

In principle, every web resource with a URI can be included and annotated in the virtual collection of our demonstrator. As a test set of data we have included three web collections:

- The Artchive collection¹² contains around 4,000 images of paintings, mainly from the 19th and 20th century.
- The ARIA collection¹³ of the Rijksmuseum in Amsterdam contains images of some 750 master pieces.
- The RMV collection¹⁴ of the Rijksmuseum voor Volkenkunde in Leiden describes about 80,000 images of ethnographic objects that belong to various cultures worldwide.

For the Artchive items we have used a parsing technique to transform the existing textual annotation in a semantic annotation, i.e. matching strings from the text to concepts from the various thesauri.

The metadata that accompagnies the Artchive collection consists of a short file holding textual values for title, creator, dimensions, material, year of creation, location and comments. Unfortunately the descriptor name is not specified with the value and not all descriptions have the same values in the same order. We used a grammar to parse and canonise the date of creation and dimension fields. Author and material are matched to ULAN and AAT using a syntactic distance measure and selecting the best match.

For the other collections we used similar strategies for enriching the original metadata with semantic categories. Adding a collection thus involves some information-extraction work on the metadata. In addition, the demonstrator supplies an manual-annotation interface which can be used to annotate any image on the Web.

¹¹ An unofficial OWL specification of the VRA elements, including links to Dublin Core, can be found at <http://e-culture.multimedien.nl/resources/>

¹² <http://www.artchive.com/>

¹³ <http://rijksmuseum.nl/aria/>

¹⁴ <http://www.rmv.nl>

7 Distributed vs. centralized collection data

The architecture is constructed to support multiple distributed image collections. Data (i.e. images) must have an external URI (we keep local copies, but that's only for caching). Ideally, we would like to get the original metadata also from external sources using standard protocols such as OAI¹⁵. In practice however, we encountered several problems with the quality of metadata retrieved via OAI, so for the moment we still depend on the local copy of the original metadata. Metadata extensions are also stored locally. In the future we hope to feed these back to the collection owners.

Vocabularies form a separate problem. The Getty vocabularies are licensed, so we cannot publish the full vocabulary as is. However, the information in the Getty vocabularies is freely accessible through the Getty online facilities¹⁶. We hope that these vocabularies will become publicly available. In the meantime, our demonstrator allows you to browse the vocabularies as a semantic structure and search for images semantically related to a vocabulary item (e.g. see Figure 2 for an example for the concept `Gold` from AAT). An RDF/OWL version of WordNet has recently been published (see above). We will move within the next months to this version (the same version as we are now using but with a different base URI).

8 Keyword search with semantic clustering

One of the goals of the demonstrator is to provide users with a familiar and simple keyword search, but still allow the user to benefit from all background knowledge from the underlying thesauri and taxonomies. The underlying search algorithm consists of several steps, that can be summarized as follows. First, it checks all RDF literals in the repository for matches on the given keyword. Second, from each match, it traverses the RDF graph until a resource of interest is found, we refer to this as a *target resource*. Finally, based on the paths from the matching literals to their target resources, the results are clustered.

To improve performance in finding the RDF literals that form the starting points, the RDF database maintains a btree index of words appearing in literals to the full literal, as well as a Porter-stem and metaphone (sounds-like) index to words. Based on these indexes, the set of literals can be searched efficiently on any logical combination of word, prefix, by-stem and by-sound matches¹⁷.

In the second step, which resources are considered of interest is currently determined by their type. The default settings return only resources of type artwork (`vra:Work`), but this can be overridden by the user. To avoid a combinatorial explosion of the search space, a number of measures had to be taken. Graph traversal is done in one direction only: always from the object in the triple to

¹⁵ <http://www.openarchives.org/>

¹⁶ See e.g. http://www.getty.edu/research/conducting_research/vocabularies/aat/ for access to the AAT.

¹⁷ See <http://www.swi-prolog.org/packages/semweb.html#sec:3.8>

the corresponding subject. Only for properties with an explicit `owl:inverseOf` relation is the graph also traversed in the other direction. While this theoretically allows the algorithm to miss out many relevant results, in practice we found that this is hardly an issue. In addition to the direction, the search space is kept under control by setting a threshold. Starting with the score of the literal match, this score is multiplied by the weight assigned to the property being traversed (all properties have been assigned a (default) weight between 0 and 1), and the search stops when the score falls under the given threshold. This approach not only improves the efficiency of the search, it also allows filtering out results with paths that are too long (which tend to be semantically so far apart, that users do not consider them relevant anymore). By setting the weights to non-default values, the search can also be fine tuned to a particular application domain.

In the final step, all results are clustered based on the path between the matching literal and the target result. When the paths are considered on the instance level, this leads to many different clusters with similar content. We found that clustering the paths on the schema level provides more meaningful results. For example, searching on keyword “fauve” matches works from Fauve painters Matisse and Derain. On the instance level, this results in different paths:

```
dc:creator -> ulan:Derain -> glink:hasStyle -> aat:fauve -> rdfs:label -> "Fauve"  
dc:creator -> ulan:Matisse -> glink:hasStyle -> aat:fauve -> rdfs:label -> "Fauve"
```

while on the schema level, this becomes a single path:

```
dc:creator -> ulan:Person -> glink:hasStyle -> aat:Concept -> rdfs:label -> "Fauve"
```

The paths are translated to English headers that mark the start of each cluster, and this already gives users an indication why the results match their keyword. The path given above results in the cluster title “Works created by an artist with matching AAT style”. To explain the exact semantic relation between the result and the keyword searched on, the instance level path is displayed when hovering over a resulting image.

9 Vocabulary and metadata statistics

Table 1 shows the number of triples that are part of the vocabularies and metadata currently being used by the demonstrator. The table has three parts: (i) the schemas (e.g. the RDF/OWL schema for WordNet defining notions such as **SynSet**), (ii) the vocabulary entries and their relationships, and (iii) the collection metadata. In total, these constitute a triple set of roughly 9,000,000 triples. We plan to extend this continuously as more collections (and corresponding vocabularies) are being added.

Acknowledgements The E-Culture project is an application subproject within the context of the Multimedien (“Multimedia Netherlands”¹⁸) project funded by the Dutch BSIK Programme.

¹⁸ <http://www.multimedien.nl>

Table 1. Number of triples for the different sources of vocabularies and collection metadata

| Document | # Sources | # Triples |
|---------------------|-----------|-----------|
| Schemas | | |
| RDFS/OWL | 2 | 358 |
| Annotation | 6 | 769 |
| Vocabularies | 9 | 1,225 |
| Collections | 1 | 29,889 |
| Vocabularies | | |
| TGN | 4 | 425,517 |
| ULAN | 16 | 1,896,936 |
| AAT | 1 | 249,162 |
| WordNet | 18 | 2,579,206 |
| Collections | | |
| Artchive | 4 | 74,414 |
| Rijkmuseum | 1 | 27,933 |
| RVM | 1 | 3,662,257 |

References

1. Wielemaker, J., Schreiber, A.T., Wielinga, B.J.: Prolog-based infrastructure for rdf: performance and scalability. In Fensel, D., Sycara, K., Mylopoulos, J., eds.: The Semantic Web - Proceedings ISWC'03, Sanibel Island, Florida. Volume 2870 of Lecture Notes in Computer Science., Berlin/Heidelberg, Springer Verlag (2003) 644–658 ISSN 0302-9743.
2. Wielemaker, J., Schreiber, G., Wielinga, B.: Using triples for implementation: the Triple20 ontology-manipulation tool. In Gil, Y., Motta, E., Benjamins, R., Musen, M., eds.: The Semantic Web – ISWC 2005: 4th International Semantic Web Conference, Galway, Ireland, November 6-10, 2005. Proceedings. Volume 3729 of Lecture Notes in Computer Science., Springer-Verlag (2005) 773–785
3. van Assem, M., Menken, M., Schreiber, G., Wielemaker, J., Wielinga, B.: A method for converting thesauri to RDF/OWL. In McLLraith, S.A., Plexousakis, D., van Harmelen, F., eds.: Proc. Third Inte. Semantic Web Conference ISWC 2004, Hiroshima, Japan. Volume 3298 of LNCS., Berlin/Heidelberg, Springer Verlag (2004) 17–31
4. van Assem, M., Gamgemi, A., Schreiber, G.: Conversion of wordnet to a standard rdf/owl representation. In: Proc. LREC 2006. (2006) Accepted for publication. <http://www.cs.vu.nl/~guus/papers/Assem06a.pdf>.
5. Miles, A., Baker, T., Swick, R.: Best practice recipes for publishing RDF vocabularies. Working draft, W3C (2006) <http://www.w3.org/TR/2006/WD-swbp-vocab-pub-20060314/>.
6. de Boer, V., van Someren, M., Wielinga, B.: Extracting instances of relations from web documents using redundancy. In: Proc. Third European Semantic Web Conference (ESWC'06), Budvar, Montenegro. (2006) Accepted for publication. <http://staff.science.uva.nl/~vdeboer/publications/eswc06paper.pdf>.
7. Visual Resources Association Standards Committee: VRA Core Categories, Version 3.0. Technical report, Visual Resources Association (2000) URL: <http://www.vraweb.org/vracore3.htm>.