

SemaPlorer—Interactive Semantic Exploration of Data and Media based on a Federated Cloud Infrastructure

Simon Schenk, Carsten Saathoff, Anton Baumesberger, Frederik Jochum,
Alexander Kleinen, Steffen Staab, and Ansgar Scherp

University of Koblenz-Landau, Germany
{sschenk,saathoff,scripper,taphyriel,alexkleinen,staab,scherp}@uni-koblenz.de
<http://isweb.uni-koblenz.de>

Abstract. SemaPlorer is an easy to use application that allows end users to interactively explore and visualize a very large, mixed-quality and semantically heterogeneous distributed semantic data set in real-time. Its purpose is to acquaint oneself about a city, touristic area, or other area of interest. By visualizing the data using a map, media, and different context views, we clearly go beyond simple storage and retrieval of large numbers of triples. The interaction with the large data set is driven by the user. SemaPlorer leverages different semantic data sources such as DBpedia, GeoNames, WordNet, and personal FOAF files. These make a significant portion of the data provided for the billion triple challenge. It intriguingly connects with a large Flickr data set converted to RDF. SemaPlorer’s storage infrastructure bases on Amazon’s Elastic Computing Cloud (EC2) and Simple Storage Service. We apply NetworkedGraphs as additional layer on top of EC2, performing as a large, federated data infrastructure for semantically heterogeneous data sources from within and outside of the cloud. Therefore, the application is scalable with respect to the amount of distributed components working together as well as the number of triples managed overall. Hence, SemaPlorer is flexible enough to leverage for exploration almost arbitrary additional data sources that might be added in the future.

1 Introduction

Informing oneself about cities, touristic regions, and other areas of interest is a task often performed on the Internet. Today’s applications supporting users in this task are centralized and monolithic such as travel sites like TripAdvisor (<http://www.tripadvisor.com>) and Wikitravel (<http://wikitravel.org>) and knowledge platforms like Freebase (<http://www.freebase.com>). With our novel infrastructure and application, SemaPlorer, we target a web of networked data spaces. Such systems, services, and data stores are easily and seamlessly integrated into a federated infrastructure in order to enable generic access to semantic multimedia data. The different data spaces may be located remotely, provided over SPARQL end points that can be queried and connected over a distributed infrastructure. (Almost) arbitrary data sources may be added ad hoc at any later point in time to extend the data infrastructure of SemaPlorer.

A major step forward towards accomplishing this visionary objective is presented in this paper with the SemaPlorer application and underlying data infrastructure. SemaPlorer is an interactive application that gives end users a usable

tool to explore and visualize a very large, mixed-quality and semantically heterogeneous distributed semantic data set in real-time. For SemaPlorer, we pursue a blended browsing and querying approach [1] to retrieve and visualize information. Users can navigate through almost arbitrary data sets using different facets (cf. [2]) such as location, time, people, and tags. When the user interacts with the application, multiple queries are sent to and executed by the underlying storage infrastructure to retrieve the appropriate results. The results are visualized using a map, media, and different context views representing the different facets.

For SemaPlorer, we have integrated and leveraged different semantic data sources such as DBpedia (<http://dbpedia.org>), GeoNames (<http://geonames.org>), WordNet (<http://wordnet.princeton.edu>), and personal FOAF files contained in the Swoogle (<http://swoogle.umbc.edu>) crawl of Semantic Web data. These make a significant portion of the data provided for the billion triple challenge. Further, we have incorporated a partial crawl of Flickr (<http://flickr.com>) as a very large non-semantic data set that has been converted to 700 million RDF triples. Together, they form a very large, semantically heterogeneous and mixed-quality data set that sums up to more than 1 billion triples. Linking this data set requires a flexible and scalable storage infrastructure. The SemaPlorer infrastructure in its current configuration consists of a set of 25 RDF stores¹. The stores are hosted on virtual machines on Amazon’s Elastic Computing Cloud (EC2, <http://aws.amazon.com/ec2/>). Amazon’s Simple Storage Service (S3, <http://aws.amazon.com/s3/>) is used to store the EC2 virtual machine images and the semantic datasets. The stores can be transparently accessed as a single, virtual RDF store through a federator. The federator uses NetworkedGraphs [4], a SPARQL-based distributed view mechanism for RDF, and distributed evaluation of SPARQL queries [5, 6]. Lightweight inferencing is done using NetworkedGraphs at runtime, e.g., for integrating semantically heterogeneous data. Thus, adding new data sources becomes extremely easy by extending the federator’s configuration while being fully transparent to the SemaPlorer application.

2 SemaPlorer Application

Collecting information about an area of interest such as a city or touristic region is a task often performed on the Internet. The more complex such queries get, the harder today’s search engines and platforms can fulfill these information requests. For example, a person interested in Berlin can easily find information about the city using standard search such as Google. However, finding places where there is some street art in the city of Berlin is almost impossible. Changing this context to another city such as Paris puts an additional challenge to the application that traditional approaches cannot solve. With the SemaPlorer application, we support the users in conducting such complex data exploration tasks. The application uses data federated from different sites using faceted, blended browsing and querying. We have defined four facets of general interest

¹ Given the scalability of today’s RDF stores, a smaller number would certainly suffice. However, this higher number illustrates the scalability of our approach with regard to federation.

in SemaPlover, namely location, time, people, and tags. Other facets can be easily configured and added. A facet provides a filtering on a large data set. For example, SemaPlover can present the sights of a certain city or area using the location facet. Blended browsing and querying means that while users interact with SemaPlover, different queries are constructed in the background and forwarded to the underlying storage infrastructure and their results are visualized on the screen. This approach allows for a user-driven visualization and interactive experience of the semantic data provided on the Web today. In SemaPlover, the users initially state a simple text query to the system as depicted in the top left corner of Fig. 1. The result list contains different places, people, and tags matching the query. When the user clicks on the city of Berlin, the SemaPlover application updates the center part of Fig. 1 showing a map of the city. Concurrently, a query is executed filling the map view with interesting places and sights, represented by pins. At the same time, a second query is executed based on what is currently seen on the map to fill the context view in the right hand side of Fig. 1.

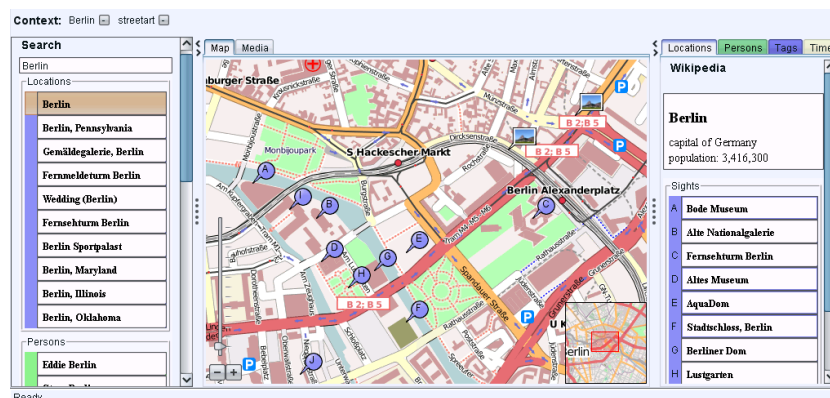


Fig. 1. Screenshot of the SemaPlover application showing street art in Berlin

For each facet, a context view is defined in the SemaPlover application. For example, the location view provides information from DBpedia such as population, country, and others. It lists sights and shows nearby places. The people view contains celebrities associated with that place, Flickr users who have uploaded geo-referenced images from that region, and Internet users living in that area according to their FOAF files. The time view allows for selecting a specific time period such as from-to-date and seasons like summer and winter. In the tag view, the tags from Flickr are shown as cloud. All elements in the context views such as sights, nearby places, celebrities, tags, and others are interactive. This means that the users can click on it to continue the blended browsing and querying. For example, when the map view shows the city of Berlin, one can click on the tag *street art*. Instantaneously, the map view is updated and locations of Flickr photos tagged as street art are shown. By stating another query for Paris, the user can switch from the current context of street art photos in Berlin and compare them with Paris.

3 SemaPlorer Dataset and Interlinking of the Data

To provide blended browsing and querying about areas of interest in SemaPlorer, different kinds of semantic data are combined. We use a significant portion of the dataset provided for the challenge, namely DBpedia (120M triples), GeoNames (70M triples), WordNet (2M triples), and Swoogle (175M triples). In addition, we use a crawl of Flickr covering several months in 2005-2006 (700M triples), which has been translated to RDF. As described in Sec. 2, we have defined different facets for our SemaPlorer application. These facets are provided by different parts of the data. In the following, we describe the data used for the different facets and how they are connected.

Location. Elements of this facet refer to geographic coordinates. We employ GeoNames for cities, countries, and others. For sights, we use a combination of full-text search on DBpedia article labels and SKOS category labels. In order to identify sights, we use the SKOS categories that are available in DBpedia. We assume that *skos:broader* is transitive and precompute the transitive closure of all resources. Subsequently, we perform a full text search on the category labels and constrain the results to resources that are connected to *dbpedia:Visitor_attractions* via *skos:subject* and the transitive closure of *skos:broader*. For displaying nearby places and sights, we select all siblings of a chosen location element and rank them based on the geo location distance. For example, when selecting the Arc de Triomphe in Paris, nearby places computed include Eiffel Tower and Notre Dame. This has to be done, because *nearbyPlace* information is missing from the GeoNames export. Images are displayed based on geo-tagged pictures on Flickr. Only locations, which are shown on flickr are provided to the user.

Time. For the time facet, there is no explicit data set defined. We rather provide the possibility to filter content from a certain time period, e.g., select pictures of a specific month from Flickr. In addition, we allow filtering of content from a particular season like winter and summer.

Person. From the datasets introduced above, we have identified three types of persons. First, we select “celebrities” from DBpedia. Second, we select users that posted images on Flickr. Finally, we search for Internet users that published their FOAF files from Swoogle. For any of these types of persons, we use a different combination of the data. For celebrities, we find images depicting the selected celebrity based on a full-text search on the Flickr tags. With respect to a Flickr user, we search for content posted by the user. For Internet users, we look at their FOAF profile’s geo location (if available) and connect it with images of that location from Flickr.

Tags. Tags are directly associated with the Flickr content. We provide full-text search over the tags. When a tag is selected by a user, we show related tags from Flickr and WordNet.

Complexity of Queries. For filling the facets described above, multiple queries are executed at the same time. For the initial search by keyword as described in Sec. 2, three simultaneous queries are performed for retrieving locations, persons, and tags. When clicking on one of the retrieved items in the search results, eight simultaneous queries are executed filling the media view and map view, calculating nearby places, selecting sights, celebrities,

Flickr users, Internet users, tags, and retrieving the DBpedia abstract. The same queries are performed when the context of the current view is changed, e.g., when the location is changed by clicking on a sight or nearby place or when a specific person or tag is selected in the corresponding context view. This approach allows for a very flexible change of the SemaPlorer application, e.g., adding certain elements to the views or removing them. The SPARQL queries make use of the full expressiveness of SPARQL, including UNION, OPTIONAL and various FILTER expressions. Additionally, Lucene queries are included in the SPARQL queries using predicate functions and Sesame LuceneSail (<http://dev.nepomuk.semanticdesktop.org/wiki/LuceneSail>). We have extended the Lucene Sail to allow for range queries and queries for geographic proximity. The queries have a standard length of 4 to 9 joins. On average, 2 to 3 joins connect multiple repositories with up to 4 datasets in a single query. As the GeoNames and Flickr datasets have been distributed over multiple repositories, a varying number of distributed unions are executed. However, these are less critical as they can easily be parallelized. Depending on the context the user selects, the queries can grow, e.g., by selecting images tagged with multiple tags in a certain time period in a certain geographic area.

Achievements and Experiences. When designing the dataset for our SemaPlorer application and working on it, we found out that the data sets are often not complete and sometimes the semantics are not explicit enough. For example, GeoNames is missing information on sights and nearby places. Nevertheless, we were able to retrieve this information by intriguingly connecting the heterogeneous data sets as described above. Considering the data set, we further observe that the data is heterogeneous even within a solitary dataset. For example, there is no clear approach for specifying the place of birth of a person in DBpedia. Sometimes it is *dbpedia:cityofbirth* and sometimes *dbpedia:birthPlace*. In SemaPlorer, we solve such ambiguities by mapping the two properties and unifying the result sets. While Linked Open Data progresses in linking the metadata, it is still an open issue how to exploit it for managing resources such as Flickr images. As SemaPlorer shows, mapping of Linked Open Data and the RDF conversion of the Flickr data is feasible and it works well, e.g., with GeoNames. However, instead of tagging images with keywords and mapping these tags with Linked Open Data, it would be more beneficial to directly use Linked Open Data to annotate the images. For example, an image depicting the Eiffel Tower could be annotated with the corresponding DBpedia concept.

4 SemaPlorer Architecture

The architecture of SemaPlorer is depicted in Fig. 2. It is divided into two subsystems: The first subsystem consists of the K-Space Annotation Tool (KAT, <https://launchpad.net/kat>) and its SemaPlorer specific extensions, the KAT Plugins. It is deployed to the client's computer and provides the user interface and application logic of the SemaPlorer application described in Sec. 2. The second subsystem implements the federated data infrastructure and comprises an Administration Component for RDF repositories, the NetworkedGraphs-based Federator, and the different RDF Stores for the semantic data and Literal Stores for the DBpedia abstracts and Flickr tags. The Administration Component and

the Federator are hosted on our local computing infrastructure. All other components, i.e., RDF Stores and Literal Stores providing the billion triple data set are hosted on Amazon EC2 nodes. The architecture of SemaPlorer and the single components are described in more detail in the following.

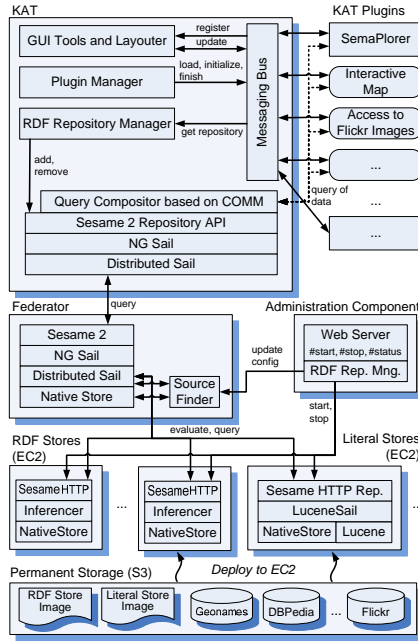


Fig. 2. Architecture of SemaPlorer

The first subsystem, provided by KAT and its plugins is a generic architecture designed to develop applications for browsing and (semi-automatically) annotating multimedia data. It can be extended by generic functionality such as an interactive map or access to Flickr images. The functionality is provided via a Messaging Bus to more application specific plugins such as the depicted SemaPlorer plugin. KAT provides a Plugin Manager for managing application specific extensions. Furthermore, it provides some GUI Tools and a GUI Layouter. Finally, KAT possesses a local storage infrastructure for multimedia annotations based on COMM [7] and Sesame 2 (<http://openrdf.org>). This storage is designed for annotations made by (semi-automatic) annotation plugins or manual annotations by the users. It will become an interesting feature for future extensions of our SemaPlorer application.

The data set described in Sec. 3 is provided through the second subsystem, the NetworkedGraphs-based federated data infrastructure leveraging Amazon’s EC2. The Administration Component of this data infrastructure controls the virtual machines running on EC2. Using a simple web GUI, EC2 nodes for specific parts of the data or the entire dataset can be started and stopped. New datasets can be created by adding a description of the dataset to a configuration file and starting the new node. Whenever nodes are started or stopped, the Administration Component updates the Federator configuration accordingly. The Federator is the single SPARQL endpoint offering SemaPlorer unified access to the whole dataset in a virtual RDF repository. Queries against the Federator are analyzed to determine, which endpoints can be used to evaluate parts of the query. Subsequently, the query is split into subqueries that are evaluated at the actual data sources [5, 6].

The dataset is stored at storage nodes in EC2 using S3. We use three different configurations for EC2 nodes: The first one stores RDF data without any inferences. It is used, e.g., for DBpedia infobox data. It also serves as basis for the other two node types. The second one uses LuceneSail and additionally provides full-text indexes over the RDF literals. It is used, e.g., for tags, DBpedia articles, and category labels. For the SemaPlorer application, we do not need full RDFS inferencing. In contrast, transitivity in SKOS hierarchies is needed, which is not

provided by RDFS. Hence, we use inferencing with custom rules in the third configuration of S3 nodes. As the custom rules inferencer of Sesame does not scale to the dataset used, we precompute the transitive closure of *skos:broader* for DBpedia categories.

In addition to SPARQL federation, the Federator performs simple schema mappings to homogenize representations from the various data sources used for SemaPlover. This schema mapping is done at run time using NetworkedGraphs. For example, for persons we have three different representations: FOAF files using the FOAF vocabulary, DBpedia using a (Living)Person category, and Flickr users. Similar challenges arise from the modeling of geographic entities and annotation of images and for providing access to properties without a clear schema such as place of birth in DBpedia. In order to allow the SemaPlover application to abstract from these differing representations, we map them to a canonical form. In the case of Persons, the FOAF vocabulary is used. As a result, we can add any dataset for which a mapping to the FOAF vocabulary is possible.

5 Related Work

The principle idea of faceted, blended browsing and querying is intriguing, but well-known, e.g., [8, 9]. The winner of the Semantic Web challenge 2006, /facet [10], has brought this idea into the arena of semantic data. Recently, the faceted application Freebase Parallax (<http://mqlx.com/~david/parallax>) emerged, a faceted browser for exploring and visualizing the structured data of Freebase (<http://www.freebase.com>). The largest disadvantage of /facet and Freebase Parallax is that they are built on a centralized infrastructure that does not allow for scalable use of a large set of data coming from many different data sources. With the SemaPlover application based on KAT and NetworkedGraphs, we have achieved this and provide for a faceted, blended browsing and querying over a very large, mixed-quality and semantically heterogeneous distributed semantic data.

Various systems providing highly scalable management of RDF data have been provided, e.g. YARS2 [11]. These systems aim at managing a large volume of RDF data in a single, albeit federated, repository. In contrast, our infrastructure aims at integrating multiple semantically heterogeneous repositories across the Semantic Web into a single virtual repository infrastructure. DARQ [12] is a related approach aiming at querying multiple SPARQL endpoints. In contrast to our system, it is based on manually maintained statistics about remote endpoints, which we do not assume to be available. Additionally, severe limitations are imposed on the structure of queries by DARQ. In the context of the Linked Open Data effort, challenges similar to our setting arise with respect to storage requirements. However, querying is not addressed. DynaQuest [13] aims at a web-scale distributed virtual relational database. However, relational databases do not cope well with semi-structured, semantically heterogeneous data.

6 Conclusions

In this paper, we have presented the SemaPlover application and data infrastructure. As shown, the SemaPlover application is an easy to use tool that allows end users to interactively explore and visualize a very large, mixed-quality distributed semantic data set in real-time. The application leverages a significant

portion of the data provided for the billion triple challenge. Further, a large Flickr data set converted to RDF is incorporated. However, the main focus of the SemaPlover application remains on the use and integration of the different data sources provided for the challenge. The storage infrastructure underlying SemaPlover allows for transparent access to arbitrary, distributed RDF repositories, in our case stored on EC2. By this, the application is scalable with respect to the amount of distributed components working together. In addition, arbitrary additional data can be added at a later point in time. Thus, using Amazon's EC2 and NetworkedGraphs brings us closer to the vision of generic access to distributed semantic multimedia data. Particularly, we have shown that besides scaling centralized repositories, connecting many smaller repositories is a feasible and in many ways a more advantageous approach to scale with regard to organizational needs of autonomous contributors on the Semantic Web.

In the long term, the preferred mode of operation will be the direct use of SPARQL endpoints run by the providers of the data. Switching to these live data sources can be easily conducted by changing the Federator's configuration and without modifying the SemaPlover application or any other application that might use the federated data infrastructure.

Acknowledgment. This research has been co-funded by the EU in FP6 in the NoE K-Space (027026) and NeOn project (027595) and FP7 in the WeKnowIt project (215453).

References

1. Munroe, K.D., Ludscher, B., Papakonstantinou, Y.: Blending Browsing and Querying of XML in a Lazy Mediator System. In: Extending Database Technology. (2000)
2. Hearst, M.A.: Design recommendations for hierarchical faceted search interfaces. In: SIGIR, Workshop on Faceted Search. (2006)
3. D. A. Smith, A. Owens, m. c. schraefel, et al.: Challenges in Supporting Faceted Semantic Browsing of Multimedia Collections. In: SAMT. (2007)
4. Schenk, S., Staab, S.: NetworkedGraphs: a declarative mechanism for SPARQL rules, SPARQL views and RDF data integration on the web. In: WWW. (2008)
5. Schenk, S., Petrak, J.: Sesame RDF Repository Extensions for Remote Querying. In: ZNALOSTI Conf. (2008)
6. Zemanek, J., Schenk, S., Svatek, V.: Optimizing sparql queries over disparate rdf data sources through distributed semi-joins. In: ISWC 2008 Poster and Demo Session Proceedings, CEUR-WS (2008)
7. Arndt, R., Troncy, R., Staab, S., Hardman, L., Vacura, M.: COMM: Designing a Well-Founded Multimedia Ontology for the Web. In: ISWC. (2007)
8. Yee, K.P., Swearingen, K., Li, K., Hearst, M.: Faceted metadata for image search and browsing. In: Human factors in computing systems, ACM (2003)
9. m. c. schraefel, D. A. Smith, A. Owens, et al.: The evolving mspace platform: leveraging the semantic web on the trail of the memex. In: Hypertext. (2005)
10. Hildebrand, M., van Ossenbruggen, J., Hardman, L.: /facet: A Browser for Heterogeneous Semantic Web Repositories. In: ISWC. (2006)
11. Harth, A., Umbrich, J., Hogan, A., Decker, S.: YARS2: A Federated Repository for Querying Graph Structured Data from the Web. In: ISWC, Springer (2007)
12. Quilitz, B., Leser, U.: Querying Distributed RDF Data Sources with SPARQL. In: ESWC. (2008)
13. Grawunder, M., Köster, F.: The DynaQuest-Framework for Dynamic and Adaptive Source Selection. In: Collaborative Technologies and Systems. (2003)