

Evaluating Semantic Search Systems to Identify Future Directions of Research

Khadija Elbedweihy¹(✉), Stuart N. Wrigley¹, Fabio Ciravegna¹,
Dorothee Reinhard², and Abraham Bernstein²

¹ University of Sheffield, Regent Court, 211 Portobello, Sheffield, UK
{k.elbedweihy,s.wrigley,f.ciravegna}@dcs.shef.ac.uk

² University of Zürich, Binzmühlestrasse 14, 8050 Zürich, Switzerland
{dreinhard,bernstein}@ifi.uzh.ch

Abstract. Recent work on searching the Semantic Web has yielded a wide range of approaches with respect to the style of input, the underlying search mechanisms and the manner in which results are presented. Each approach has an impact upon the quality of the information retrieved and the user's experience of the search process. This highlights the need for formalised and consistent evaluation to benchmark the coverage, applicability and usability of existing tools and provide indications of future directions for advancement of the state-of-the-art. In this paper, we describe a comprehensive evaluation methodology which addresses both the underlying performance and the subjective usability of a tool. We present the key outcomes of a recently completed international evaluation campaign which adopted this approach and thus identify a number of new requirements for semantic search tools from both the perspective of the underlying technology as well as the user experience.

1 Introduction and Related Work

State-of-the-art semantic search approaches are characterised by their high level of diversity both in their features as well as their capabilities. Such approaches employ different styles for accepting the user query (e.g., forms, graphs, keywords) and apply a range of different strategies during processing and execution of the queries. They also differ in the format and content of the results presented to the user. All of these factors influence the user's perceived performance and usability of the tool. This highlights the need for a formalised and consistent evaluation which is capable of dealing with this diversity. It is essential that we do not forget that searching is a user-centric process and that the evaluation mechanism should capture the usability of a particular approach.

One of the very first evaluation efforts in the field was conducted by Kaufmann to compare four different query interfaces [1]. Three were based on natural language input (with one employing a restricted query formulation grammar); the fourth employed a formal query approach which was hidden from the

This work was supported by the European Union 7th FP ICT based e-Infrastructures Project SEALS (Semantic Evaluation at Large Scale, FP7-238975).

end user by a graphical query interface. Recently, evaluating semantic search approaches gained more attention both in IR – within its most established evaluation conference TREC – [2] as well as in the Semantic Web community (Sem-Search [3] and QALD¹ challenges).

The above evaluations are all based upon the Cranfield methodology [4]²: using a test collection, a set of tasks and a set of relevance judgments. This leaves aside aspects of user-oriented evaluations concerned with the usability of the evaluated systems and the user experience which is as important as assessing the performance of the systems. Additionally, the above attempts are separate efforts lacking standardised evaluation approaches and measures. Indeed, Halpin et al. [3] note that “the lack of standardised evaluation has become a serious bottleneck to further progress in this field”.

The first part of this paper describes an evaluation methodology for assessing and comparing the strengths and weaknesses of user-focussed Semantic Search approaches. We describe the dataset and questions used in the evaluation and discuss the results of the usability study. The analysis and feedback from this evaluation are described. The second part of the paper identifies a number of new requirements for search approaches based upon the outcomes of the evaluation and analysis of the current state-of-the-art.

2 Evaluation Design

In the Semantic Web community, semantic search is widely used to refer to a number of different categories of systems:

- *gateways* (e.g., Sindice [5] and Watson [6]) locating ontologies and documents
- approaches reasoning over data and information located within documents and ontologies (PowerAqua [7] and FREyA [8])
- view-based interfaces allowing users to explore the search space while formulating their queries (Semantic Crystal [9], K-Search [10] and Smeagol [11])
- mashups integrating data from different sources to provide rich descriptions about Semantic Web objects (Sig.ma [12]).

The evaluation described here focuses on user-centric semantic search tools (e.g. query given as keywords or natural language or using a form or a graph) querying a repository of semantic data and returning answers extracted from them. The tools’ results presentation is not limited to a specific style (e.g., list of entity URIs or a visualisation of the results). However, the results returned must be answers rather than documents matching the given query.

Search is a user-centric activity; therefore, it is important to emphasise the users’ experience. An important aspect of this is the formal gathering of feedback from the participants which should be achieved using standard questionnaires.

¹ <http://www.sc.cit-ec.uni-bielefeld.de/qald-1>.

² <http://www.sigir.org/museum/pdfs/ASLIB%20CRANFIELD%20RESEARCH%20PROJECT-1960/pdfs/>.

Furthermore, the use of an additional demographics questionnaire allows more in-depth findings to be identified (e.g., if a particular type of user prefers a particular search approach).

2.1 Datasets and Questions

Subjects are asked to reformulate a set of questions using a tool’s interface. Thus, it is important that the data set would be from an understandable and well-known domain (and hence, easily understandable by non-expert users) and, preferably, already have a set of questions and associated groundtruths. The geographical dataset from the Mooney Natural Language Learning Data³ satisfies these requirements and has been used in a number of usability studies [1,8]. Although the Mooney dataset is different from ones currently found on the Semantic Web such as DBpedia in terms of size, heterogeneity and quality, the assessment of the tools ability to handle these aspects is not the focus of this phase but rather the usability of the tools and the user experience.

The questions [13] used in the first evaluation campaign were generated based on the existing templates within the Mooney dataset. These contained questions with varying complexity and assessing different features. For instance, they contained simple with only 1 unknown concept such as “Give me all the capitals of the USA?” and comparative questions such as “Which rivers in Arkansas are longer than Aleghany river” as well as negation questions such as “Tell me which river do not traverse the state with capital nashville”.

2.2 Criteria and Analyses

Usability. Different input styles (e.g., form-based, NL, etc.) can be compared with respect to the input query language’s expressiveness and usability. These concepts are assessed by capturing feedback regarding the user experience and the usefulness of the query language in supporting users to express their information needs and formulate searches [14]. Additionally, the expressive power of a query language specifies what queries a user is able to pose [15]. The usability is further assessed with respect to results presentation and suitability of the returned answers (data) to the casual users as perceived by them. The datasets and associated questions were designed to fully investigate these issues.

Performance. Users are familiar with the performance of commercial search engines (e.g., Google) in which results are returned within fractions of a second; therefore, it is a core criterion to measure the tool’s performance with respect to the speed of execution.

Analyses. The experiment was controlled using custom-written software which allowed each experiment run to be orchestrated and timings and results to be captured. The results included the actual result set returned by a tool for a

³ <http://www.cs.utexas.edu/users/ml/nldata.html>.

query, the time required to execute a query, the number of attempts required by a user to obtain a satisfactory answer as well as the time required to formulate the query. We used post-search questionnaires to collect data regarding the user experience and satisfaction with the tool. Three different types of online questionnaires were used which serve different purposes. The first is the System Usability Scale (SUS) questionnaire [16]. The test consists of ten normalized questions and covers a variety of usability aspects, such as the need for support, training, and complexity and has proven to be very useful when investigating interface usability [17]. We developed a second, extended, questionnaire which includes further questions regarding the satisfaction of the users. This encompasses the design of the tool, the input query language, the tool's feedback, and the user's emotional state during the work with the tool. An example of a question used is *'The query language was easy to understand and use'* with answers represented on a scale from 'disagree' to 'agree'. Finally, a demographics questionnaire collected information regarding the participants.

3 Evaluation Execution and Results

The evaluation consisted of tools from form-based, controlled-NL-based and free-NL-based approaches. Each tool was evaluated with 10 subjects (except K-Search [10] which had 8) totalling 38 subjects (26 males, 12 females) aged between 20 and 35 years old. They consisted of 28 students and 10 researchers drawn from the University population. Subjects rated their knowledge of the Semantic Web with 6 reporting their knowledge to be advanced, 5 good, 9 average, 10 little and 8 having no experience. In addition, their knowledge of query languages was recorded, with 5 stating their knowledge to be advanced, 12 good, 8 average, 6 little and 7 having no experience.

Firstly, the subjects were presented with a short introduction to the experiment itself such as why the experiment is taking place, what is being tested, how the experiment will be executed, etc. Then the tool itself was explained to the subjects; they learnt about the type and the functionality of the tool and how to apply its specific query language to answer the given tasks. The users were then given sample tasks to test their understanding of the previous phases. After that, the subjects did the actual experiment: using the tool's interface to formulate each question and get the answers. Having finished all the questions, they were presented with the three questionnaires (Sect. 3). Finally, the subjects had the chance to talk about important and open questions and give more feedback and input to their satisfaction or problems with the system being tested.

Table 1 shows the results for the four tools participating in this phase. The *mean number of attempts* shows how many times the user had to reformulate their query in order to obtain answers with which they were satisfied (or indicated that they were confident a suitable answer could not be found). This latter distinction between finding the appropriate answer and the user 'giving up' after a number of attempts is shown by the *mean answer found rate*. *Input time* refers to the amount of time the subject spent formulating their query using the tool's interface, which acts as a core indicator of the tool's usability.

Table 1. Evaluation results showing the tools performance. Rows refer to particular metrics.

Criterion	K-Search <i>Form-based</i>	Ginseng <i>Controlled NL-based</i>	NLP-Reduce <i>NL-based</i>	PowerAqua <i>NL-based</i>
Mean experiment time (s)	4313.84	3612.12	4798.58	2003.9
Mean SUS (%)	44.38	40	25.94	72.25
Mean ext. questionnaire (%)	47.29	45	44.63	80.67
Mean number of attempts	2.37	2.03	5.54	2.01
Mean answer found rate	0.41	0.19	0.21	0.55
Mean execution time (s)	0.44	0.51	0.51	11
Mean input time (s)	69.11	81.63	29	16.03

According to the ratings of SUS scores [18], none of the four participating tools fell in either the best or worst category. Only one of the tools (PowerAqua [7]) had a ‘Good’ rating with a SUS score of 72.25, other two tools (Ginseng [19] and K-Search [10]) fell in the ‘Poor’ rating while the last one (NLP-Reduce [20]) was classified as ‘Awful’. The results of the questionnaires were confirmed by the recorded usability measures. Subjects using the tool with the lowest SUS score (NLP-Reduce) required more than twice the number of attempts of the other tools before they were satisfied with the answer or moved on. Similarly, subjects using the two tools with the highest SUS and extended scores (PowerAqua and K-Search) found satisfying answers to their queries twice the times as for the other tools. Altogether, this confirms the reliability of the results and the feedback of the users and also the conclusions based on them.

4 Usability Feedback and Analysis

This section discusses the results and feedback collected from the subjects of the usability study. Figure 1 summarises the features most liked and disliked based on their feedback. The following discussion stems from this summary.

4.1 Input Style

On the one hand, Uren et al. [14] state that forms can be helpful to explore the search space when it is unknown to the users. Additionally, Corese [21] – which uses a form-based interface to allow users to build their queries – received very positive comments from its users among which was an appreciation for its form-based interface. On the other hand, Lei et al. [22] see this exploration as a burden on users that requires them to be (or become) familiar with the underlying ontology and semantic data. The results of our evaluation and the feedback from the users support both arguments: positive comments for the form-based

tool (K-Search) included ones such as “I liked to see the concepts and relations between them, it helped in knowing what sort of information is available to be retrieved from the system”. On the other hand, negative comments included ones such as “For me, it was complex to build some queries” and “It was hard to understand without explanation and it seemed less intuitive than NL-based tools”.

Additionally, we found that form-based interfaces allow users to build more complex queries than the natural language interfaces. However, building queries by exploring the search space is usually time consuming especially as the ontology gets larger or the query gets more complex. This was shown by Kaufmann et al. [1] in their usability study which found that users spent the most time when working with the graph-based system *Semantic Crystal*. Our evaluation supports this general conclusion: subjects using the form-based approach took between two to three times the time taken by users of natural language approaches. Also, feedback showed that most of the users found query formulation with the form-based tool (K-Search) to be laborious and requiring long time especially when they compared it to the NL-based tools (NLP-Reduce and PowerAqua). However, our analysis suggests a more nuanced behaviour. While freeform natural language interfaces are generally faster in terms of query formulation, we found this did not hold for approaches employing a very restricted language model. For instance, query formulation took longer using Ginseng (restricted natural language) than K-Search (form-based). This is further supported by users feedback: the most repeated positive comment for the free-NL-based tools was “It is quick and easy to use”. On the other hand, the more time required to formulate queries in Ginseng was due to its restrictive model which limited users expressivity and affected their satisfaction. Some of the most negative repeated comments for Ginseng included:

- It was unclear how individual terms suggested by the tool related to particular classes or relations.
- In most of the queries, I got stuck and could no longer complete the query in the way I wanted because it was restricted.
- It was very annoying and frustrating.

Kaufmann et al. [1] also showed that a natural language interface was judged by users to be the most useful and best liked. Their conclusion, that this was because users can communicate their information needs far more effectively when using a familiar and natural input style, is supported by our findings. The same study found that people can express more semantics when they use full sentences as opposed to simply keywords. Similarly, Demidova et al. [23] state that natural language queries offer users more expressivity to describe their information needs than keywords – a finding also confirmed by the user feedback from our study.

However, natural language approaches suffer from both syntactic as well as semantic ambiguities. This makes the overall performance of such approaches heavily dependent upon the performance of the underlying natural language processing techniques responsible for parsing and analysing the users’ natural

	Liked/Required	Disliked
Input Style	<ul style="list-style-type: none"> View search domain Build complex queries (AND, OR,...) Auto-completion Easy & fast input Natural & familiar language 	<ul style="list-style-type: none"> Input format complexity Restricted language model Requires knowledge of ontologies No support for superlatives or comparatives in queries Abstraction of search domain
Query Execution	<ul style="list-style-type: none"> Feedback during query execution 	<ul style="list-style-type: none"> Slow response No incremental results
Results Presentation	<ul style="list-style-type: none"> Merging results Show provenance of results 	<ul style="list-style-type: none"> Not suitable for casual users No storing/reuse of query results No sorting, grouping, or filtering of results

Fig. 1. Summary of evaluation feedback: features most liked and disliked by users categorised with respect to query format, query execution, and results presentation.

language sentences. This was shown by the feedback we received from users of the NL-based tool (NLP-Reduce), one of which was “the response is very dependent on the use of the correct terms in the query”. This was also confirmed by that approach achieving the lowest precision. Another limitation faced by the natural language approach is the lack of knowledge of the underlying ontology terms and relations by the users due to the high abstraction of the search domain. The effect of this is that any keywords or terms used by users are likely to be very different from the semantically-corresponding terms in the ontology. This in turn increases the difficulty of parsing the user query and affects the performance.

Using a restricted grammar as employed by Ginseng is an approach to limit the impact of both of these problems. The ‘autocompletion’ provided by the system based on the underlying grammar attempts to bridge the domain abstraction gap and also resembles the form-based approach in helping the user to better understand the search space. Although it provides the user with knowledge regarding which concepts, relations and instances are found in the search space and hence can be used to build valid queries, it still lacks the power of visualising the structure of the used ontology. The impact of this ‘intermediate’ functionality can be observed in the users feedback with a lower degree of dissatisfaction regarding the ability to conceptualise the underlying data but still not completely eliminated. For instance, the positive comment “The autocompletion and suggestions was helpful to know the underlying data” was given by some of the users for Ginseng. The restricted language model also prevents unacceptable/invalid queries in the used grammar by employing a guided input natural language approach. However, only accepting specific concepts and relations – found in the grammar – limits the flexibility and expressiveness of the

user queries. User coercion into following predefined sentence structures proves to be frustrating and too complicated [1, 24]. Again, this was supported by the feedback showing that users were often annoyed by this restriction especially when they got stuck and did not know how to continue the query formulation.

The feedback from the questionnaires showed that using superlatives or comparatives in the user queries (e.g.: highest point, longer than) was not supported by any of the participating tools; an issue raised by 8 subjects in the answer of the SUS question “What didn’t you like about the system and why?” and by others in the open feedback after the experiment. Only one provided a feature similar to this functionality: the ability to specify a range of values for numeric datatypes. A comparative such as *less than 5000* could then be translated to the range *0 to 5000*. However, this was deemed to be both confusing (since the user had to decide what to use as the non-specified bound) and, when the non-specified bounds were incorrect, having a negative impact on the results.

4.2 Query Execution and Response Time

Speed of response is an important factor for users since they are used to the performance of commercial search engines (e.g., Google) in which results are returned within fractions of a second. Many users in our study were expecting similar performance from the semantic search tools. Although the average response time of three of the tools (K-Search, NLP-Reduce, Ginseng) is less than a second (44 ms, 51 ms, and 51 ms respectively), users reported their dissatisfaction with these timings especially the ones who evaluated PowerAqua with response time of 11 s on average. The lack of feedback on the status of the execution process only served to increase the sense of dissatisfaction: no tool indicated the execution progress or whether a problem had occurred in the system. This lack of feedback resulted in users suspecting that something had gone wrong with the system – even if the search was still progressing – and start a new search. Furthermore, some tools made it impossible to distinguish between an empty result set, a problem with the query formulation or a problem with the search. This not only affected the users experience and satisfaction but also the approach’s measured performance. The following list includes some of the most repeated comments given by the users for all the tools in this context:

- With some queries, I had no idea whether the search was in progress or failed since there was no feedback about it.
- When the response was delayed, I suspected that an error occurred and I restarted the search process.
- It was confusing when I got back no results and I didn’t know whether this was due to an error in the system or my query, or there was actually no results for my question.

4.3 Results Presentation

Semantic Search tools are different from Semantic Web gateways or entry points such as Watson and Sindice. The latter are not intended for casual users but

for other applications or the Semantic Web community to locate Semantic Web resources such as ontologies or Semantic Web documents and are usually presented as a set of URIs. For example, Sindice shows the URIs of documents and, for every document, it additionally presents the triples contained within the document, an RDF graph of the triples, and the used ontologies. Semantic Search tools are, on the other hand, used by casual users (i.e., users who may be experts in the domain of the underlying data but may have no knowledge of semantic technologies). Such users usually have different requirements and expectations of *what* and *how* results should be presented to them.

In contrast to these ‘casual user’ requirements, a number of the search tools did not present their results in a user-friendly manner and this was reflected in the feedback. Two approaches presented the full URIs together with the concepts in the ontology that were matched with the terms in the user query. Another used the instance labels to provide a natural language presentation; however, such labels (e.g., ‘montgomeryAl’) were not necessarily suitable for direct inclusion into a natural language phrase. Indeed, the tool also displayed the ontologies used as well as the mappings that were found between the ontology and the query terms. Although potentially useful to an expert in the semantic web field, this was not helpful to casual users. In this context, some of the negative comments repeated for most of the tools include:

- I found the URIs and sometimes ontology triples presented were technical and more targeted to experts.
- It would be good to change the way results are presented to allow non-experts to understand it.

The other commonly reported limitation of all the tools was the degree to which a query’s results could be stored or reused. A number of the questions used in the evaluation had a high complexity level and needed to be split into two or more sub-queries. For instance, for the question “Which rivers in Arkansas are longer than the Allegheny river?”, the users were first querying the data for the length of the Allegheny river and then performing a second query to find the rivers in Arkansas which are longer than the answer they got. Therefore, users often wanted to use previous results as the basis of a further query or to temporarily store the results in order to perform an intersection or union operation with the current result set. Unfortunately, this was not supported by any of the participating tools. However, this shows that users have very high expectations of the usability and functionalities offered by a semantic search tool as this requirement is not provided even by traditional search systems (e.g., Google and Yahoo). Another means of managing the results that users requested was the ability to filter results according to some suitable criteria and checking the provenance of the results; only one tool provided the latter. Indeed, even basic manipulations such as sorting were requested – a feature of particular importance for tools which did not allow query formulations to include superlatives. Again, some of the comments repeated in this context included:

- It would be nice to be able to sort the results.
- I often wanted to store previous answers and use them in subsequent queries or merge them with future result sets.

5 Future Directions

This section identifies a number of areas for improvement for semantic search tools from the perspective of the underlying technology and the user experience. It is motivated by the findings previously discussed in Sect. 4 and thus a similar structure is used.

5.1 Input Style

Usability. The feedback shows that it’s very helpful for users – especially those who are unfamiliar with the underlying data – to explore the search space while building their queries using view-based interfaces which expose the structure of the ontology in a graphical manner. It gives users a much better understanding of what information is available and what queries are supported by the tool. In contrast, the feedback also shows that, when creating their queries, users prefer natural language interfaces because they are quick and easy. Clearly both approaches have their advantages; however, they suffer from various limitations when used separately as discussed in Sect. 4.1. Therefore, we believe that the combination of both approaches would help get the best of both worlds.

Users not familiar with the search domain can use a form-based or natural language-based interface to build their queries. Simultaneously, the tool should dynamically generate a visual representation of the user’s query based upon the structure of the ontology. Indeed, the user should be able to move from one query formulation style to another – at will – with each being updated to reflect changes made in the other. This ‘dual’ query formulation would ensure a casual user correctly formulates their intended query. Expert users, or those who find it laborious to use the visual approach, would simply use the natural language input facility provided by the tool. The visualisation of the query structure is still useful for such users. An additional feature for natural language input would be an *optional* ‘auto-completion’ feature which could guide the user to query completion given knowledge of the underlying ontology. This not only helps the user but also alleviates the problem of mismatching the user terms with the correct ones in the underlying search space.

Expressiveness. The feedback also shows that the evaluated tools had difficulties with supporting complex queries such as the ones containing logical operators (e.g., “AND”). Allowing the user to input more than one query and combining them with their chosen logical operator from a list included in the interface would reduce the impact of this limitation. The tool would merge the results according to the used operator (e.g., “intersection” for “AND”). For instance, a query such as “What are the rivers that pass through California and

Query	Suggestions
Which city has the largest population in California?	<ol style="list-style-type: none"> 1. Max (city population) 2. Min (city population) 3. Sum (city population) 4. None

Fig. 2. Suggestions generated by FREyA [8] for a datatype property, to handle superlatives and comparatives.

Arizona?” would be constructed as two subqueries: “What are the rivers that pass through California?” and “What are the rivers that pass through Arizona?” with the final results being the intersection of both result sets.

Furthermore, the evaluated tools faced similar difficulties with supporting superlatives and comparatives in users’ queries. FREyA [8] deals with this problem by asking the user to identify the correct choice from a list of suggestions. To illustrate this we’ll use the query “Which city has the largest population in California?”. If the system captures a concept in the user query that is a datatype property of type number, it generates maximum, minimum and sum functions. The generated suggestions for our query example are shown in Fig. 2. The user can then choose the correct superlative or comparative depending on their needs. A similar approach can be used to allow the use of superlatives and comparatives in natural language interfaces and form-based interface. In the case of the latter, whenever a datatype property is selected by the user, the tool can allow them to select from a list of functions that cover superlatives and comparatives (e.g., ‘maximum’, ‘minimum’, ‘more than’, ‘less than’).

5.2 Query Execution and Response Time

Several users reported dissatisfaction with the tools’ response time to some of their queries. Users appreciated the fact that the tools returned more accurate answers than they would get from traditional search engines, however this did not remove the effect of the delay in response – even if it was relatively small. Additionally, the study found that the use of feedback reduces the effect of the delay; users showed greater willingness to wait if they were informed that the search is still being performed and that the delay is not due to a failure in the system.

The presentation of intermediate, or partially complete, results reduces the perceived delay associated with the complete result set (e.g., Sig.ma [12]). Although only partial results are available initially, it provides both feedback that the search is executing properly and allows the user to start thinking about the content of the results before the complete set is ready. However, it ought to be noted that this approach may induce confusion in the user as the screen content changes rapidly for a number of seconds. Adequate feedback is essential even for tools which exhibit high performance and good response times. Delays may occur at a number of points in the search process and may be the result of influences beyond the developer’s control (e.g., network communication delays).

5.3 Results Presentation

Most of the users were frustrated by the fact that they didn't understand the results presented to them, feeling that too much technical knowledge was assumed. The evaluation showed that the tools underestimated the effect of this on user's experience and satisfaction.

Query answers ought to be presented to users in an accessible and attractive manner. Indeed, the tool should go a step further and augment the direct answer with associated information in order to provide a 'richer' experience for the user. This approach is adopted by WolframAlpha⁴. For example, as shown in Fig. 3, the response to the query 'What is the capital of Alabama?' includes the natural language presentation of the answer (A) as well as various population statistics (B), a map showing the location of the city (C), and other related information such as the current local time, weather and nearby cities.

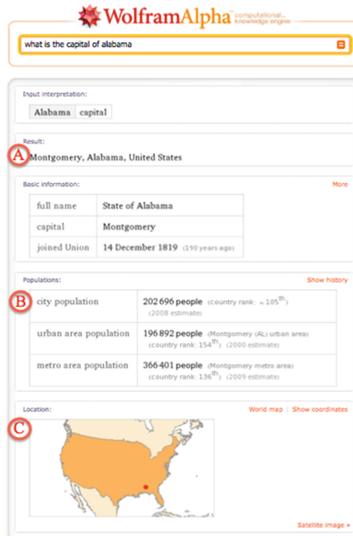


Fig. 3. Results presentation in *WolframAlpha*

An interesting requirement found by our study was the ability to store the result set of a query to use in subsequent queries. This would allow more complex questions to be answered which, in turn, improves the tools' performance. QuiKey [25] provides a functionality similar to this. QuiKey is an interaction approach that offers interactive fine grained access to structured information sources in a lightweight user interface. It allows a query to be saved which can later be used for building other queries. More complex queries can be constructed by combining saved queries with logical operators such as 'AND' and 'OR'.

⁴ <http://www.wolframalpha.com/>.

Fig. 4. Example of a Sig.ma profile

Result management was also identified as being of importance to users with commonly requested functionality included sorting, filtering and more complex activities such as establishing the provenance and trustworthiness of certain results. For example, Sig.ma [12] – a system built on top of Sindice⁵ – creates information aggregates called Entity Profiles and provides users with various capabilities to organise, use and establish the provenance of the results. Figure 4 shows part of the sigma for ‘Fabio Ciravegna’. As shown in the figure, users can see all the sources contributing to a specific profile (A) and approve or reject certain ones (B), thus filtering the results. They can also check which values in the profile are given by a specific source as they get highlighted whenever the user scrolls over the source (C) thus checking provenance of the results. Sig.ma also supports the aspect of merging separate results by allowing users to view ones returned only from selected sources.

6 Conclusions

We have presented a flexible and comprehensive methodology for evaluating different semantic search approaches; we have also highlighted a number of empirical findings from an international semantic search evaluation campaign based upon this methodology. Finally, based upon analysis of the evaluation outcomes, we have described a number of additional requirements for current and future semantic search solutions.

In contrast to other benchmarking efforts, we emphasised the need for an evaluation methodology which addressed both performance and usability [24]. We presented the core criteria that must be evaluated together with a discussion of the main outcomes. This analysis identified two core findings which impact upon semantic search tool requirements.

⁵ <http://sindice.com/>.

Firstly, we found that an intelligent combination of natural language and view-based input styles would provide a significant increase in search effectiveness and user satisfaction. Such a ‘dual’ query formulation approach would combine the ease with which a view-based approach can be used to explore and learn the structure of the underlying data whilst still being able to exploit the efficiency and simplicity of a natural language interface.

Secondly, (and perhaps of greatest interest to users) was the need for more sophisticated results presentation and management. Not only should fine grain feedback be provided on the progress of the search (such as the provision of intermediate results) but the final results should allow a large degree of customisability (sorting, filtering, saving of intermediate results, augmenting, etc.). Indeed, it would also be beneficial to provide data which is supplementary to the original query to increase ‘richness’. Furthermore, users expect to be able to have immediate access to provenance information.

In summary, this paper has presented a number of important findings which are of interest both to semantic search tool developers but also designers of interactive search evaluations. Such evaluations (and the associated analyses as presented here) provide the impetus to improve search solutions and enhance the user experience.

References

1. Kaufmann, E.: Talking to the Semantic Web – Natural Language Query Interfaces for Casual End-Users. Ph.D. thesis, University of Zurich (2007)
2. Balog, K., Serdyukov, P., de Vries, A.P.: Overview of the TREC 2011 entity track. In: TREC 2011 Working Notes (2011)
3. Halpin, H., Herzig, D.M., Mika, P., Blanco, R., Pound, J., Thompson, H.S., Tran, D.T.: Evaluating Ad-Hoc object retrieval. In: Proceedings of the IWEST 2010 Workshop (2010)
4. Cleverdon, C.W.: Report on the first stage of an investigation onto the comparative efficiency of indexing systems. Technical report, The College of Aeronautics, Cranfield, England (1960)
5. Tummarello, G., Delbru, R., Oren, E.: Sindice.com: weaving the open linked data. In: Aberer, K., et al. (eds.) ASWC 2007 and ISWC 2007. LNCS, vol. 4825, pp. 552–565. Springer, Heidelberg (2007)
6. d’Aquin, M., Baldassarre, C., Gridinoc, L., Angeletou, S., Sabou, M., Motta, E.: Characterizing knowledge on the semantic web with watson. In: EON, pp. 1–10 (2007)
7. Lopez, V., Motta, E., Uren, V.S.: PowerAqua: fishing the semantic web. In: Sure, Y., Domingue, J. (eds.) ESWC 2006. LNCS, vol. 4011, pp. 393–410. Springer, Heidelberg (2006)
8. Damljanovic, D., Agatonovic, M., Cunningham, H.: Natural language interfaces to ontologies: combining syntactic analysis and ontology-based lookup through the user interaction. In: Aroyo, L., Antoniou, G., Hyvönen, E., ten Teije, A., Stuckenschmidt, H., Cabral, L., Tudorache, T. (eds.) ESWC 2010, Part I. LNCS, vol. 6088, pp. 106–120. Springer, Heidelberg (2010)

9. Bernstein, A., Kaufmann, E., Göhring, A., Kiefer, C.: Querying ontologies: a controlled english interface for end-users. In: Gil, Y., Motta, E., Benjamins, V.R., Musen, M.A. (eds.) ISWC 2005. LNCS, vol. 3729, pp. 112–126. Springer, Heidelberg (2005)
10. Bhagdev, R., Chapman, S., Ciravegna, F., Lanfranchi, V., Petrelli, D.: Hybrid search: effectively combining keywords and semantic searches. In: Bechhofer, S., Hauswirth, M., Hoffmann, J., Koubarakis, M. (eds.) ESWC 2008. LNCS, vol. 5021, pp. 554–568. Springer, Heidelberg (2008)
11. Clemmer, A., Davies, S.: Smeagol: a “specific-to-general” semantic web query interface paradigm for novices. In: Hameurlain, A., Liddle, S.W., Schewe, K.-D., Zhou, X. (eds.) DEXA 2011, Part I. LNCS, vol. 6860, pp. 288–302. Springer, Heidelberg (2011)
12. Tummarello, G., Cyganiak, R., Catasta, M., Danielczyk, S., Delbru, R., Decker, S.: Sig.ma: live views on the web of data. In: Proceedings of the WWW 2010 (2010)
13. Wrigley, S.N., Elbedweihy, K., Reinhard, D., Bernstein, A., Ciravegna, F.: D13.3 Results of the first evaluation of semantic search tools. Technical report, SEALS Consortium (2010)
14. Uren, V., Lei, Y., Lopez, V., Liu, H., Motta, E., Giordanino, M.: The usability of semantic search tools: a review. *Knowl. Eng. Rev.* **22**, 361–377 (2007)
15. Angles, R., Gutierrez, C.: The expressive power of SPARQL. In: Sheth, A.P., Staab, S., Dean, M., Paolucci, M., Maynard, D., Finin, T., Thirunarayan, K. (eds.) ISWC 2008. LNCS, vol. 5318, pp. 114–129. Springer, Heidelberg (2008)
16. Brooke, J.: SUS: a quick and dirty usability scale. In: Usability Evaluation in Industry, pp. 189–194. CRC Press (1996)
17. Bangor, A., Kortum, P.T., Miller, J.T.: An empirical evaluation of the system usability scale. *Int. J. Hum. Comput. Interact.* **24**(6), 574–594 (2008)
18. Bangor, A., Kortum, P.T., Miller, J.T.: Determining what individual SUS scores mean: adding an adjective rating scale. *J. Usability Stud.* **4**(3), 114–123 (2009)
19. Bernstein, A., Kaufmann, E., Kaiser, C.: Querying the semantic web with ginseng: a guided input natural language search engine. In: Proceedings of the WITS 2005 Workshop (2005)
20. Kaufmann, E., Bernstein, A., Fischer, L.: NLP-reduce: a “naïve” but domain-independent natural language interface for querying ontologies. In: Proceedings of the ESWC 2007 (2007)
21. Corby, O., Dieng-Kuntz, R., Faron-Zucker, C., Gandon, F.: Searching the semantic web: approximate query processing based on ontologies. *IEEE Intell. Syst.* **21**, 20–27 (2006)
22. Lei, Y., Uren, V.S., Motta, E.: SemSearch: a search engine for the semantic web. In: Staab, S., Svátek, V. (eds.) EKAW 2006. LNCS (LNAI), vol. 4248, pp. 238–245. Springer, Heidelberg (2006)
23. Demidova, E., Nejdl, W.: Usability and expressiveness in database keyword search: bridging the gap. In: Proceedings of the VLDB Ph.D. Workshop (2009)
24. Wrigley, S.N., Elbedweihy, K., Reinhard, D., Bernstein, A., Ciravegna, F.: Evaluating semantic search tools using the SEALS platform. In: Proceedings of the IWEST 2010 Workshop (2010)
25. Haller, H.: QuiKey – an efficient semantic command line. In: Cimiano, P., Pinto, H.S. (eds.) EKAW 2010. LNCS, vol. 6317, pp. 473–482. Springer, Heidelberg (2010)