

12-2003

# Semantic (Web) Technology in Action: Ontology Driven Information Systems for Search, Integration, and Analysis


Amit P. Sheth

Wright State University - Main Campus, amit.sheth@wright.edu

Cartic Ramakrishnan

Wright State University - Main Campus

Follow this and additional works at: <http://corescholar.libraries.wright.edu/knoesis>

 Part of the [Bioinformatics Commons](#), [Communication Technology and New Media Commons](#), [Databases and Information Systems Commons](#), [OS and Networks Commons](#), and the [Science and Technology Studies Commons](#)

## Repository Citation

Sheth, A. P., & Ramakrishnan, C. (2003). Semantic (Web) Technology in Action: Ontology Driven Information Systems for Search, Integration, and Analysis. *IEEE Data Engineering Bulletin*, 26 (4), 40-48.  
<http://corescholar.libraries.wright.edu/knoesis/970>

This Article is brought to you for free and open access by the The Ohio Center of Excellence in Knowledge-Enabled Computing (Kno.e.sis) at CORE Scholar. It has been accepted for inclusion in Kno.e.sis Publications by an authorized administrator of CORE Scholar. For more information, please contact [corescholar@www.libraries.wright.edu](mailto:corescholar@www.libraries.wright.edu).

## Semantic (Web) Technology In Action: Ontology Driven Information Systems for Search, Integration and Analysis

Amit Sheth<sup>1,2</sup> and Cartic Ramakrishnan<sup>2</sup>  
<sup>1</sup>Semagix and <sup>2</sup>LSDIS lab, University of Georgia

### Abstract

Semantics is seen as the key ingredient in the next phase of the Web infrastructure as well as the next generation of information systems applications. In this context, we review some of the reservations expressed about the viability of the Semantic Web. We respond to these by identifying a Semantic Technology that supports the key capabilities also needed to realize the Semantic Web vision, namely representing, acquiring and utilizing knowledge. Given that scalability is a key challenge, we briefly review our observations from developing three classes of real world applications and corresponding technology components: search/browsing, integration, and analytics. We distinguish this proven technology from some parts of the Semantic Web approach and offer subjective remarks which we hope will foster additional debate.

### 1. Introduction

Semantics is arguably the single most important ingredient in propelling the Web to its next phase, and is closely supported by Web services and Web processes that provide standards based interoperability of applications. Semantics is considered to be the best framework to deal with the heterogeneity, massive scale, and dynamic nature of the resources on the Web. Issues pertaining to semantics have been addressed in other fields like linguistics, knowledge representation, and AI. The promise of semantics and challenges in developing semantic techniques are not new to researchers in the database and information system field either. For instance, semantics has been studied or applied in the context of data modeling, query and transaction processing, etc. Recently, a group of both database and non-database researchers came together at the *Amicalola State Park* for an intensive look at the relationship between database research and the Semantic Web. During this collaboration, they identified three pages worth of opportunities to further database research while addressing the challenges in realizing the Semantic Web [Sheth and Meersman 2002]. A follow on workshop also presented opportunity to present research at the intersection of database and the Semantic Web [<http://swdb.semanticweb.org>].

Nevertheless, many researchers in the database community continue to express significant reservations toward the Semantic Web. **Table 1** shows some examples of criticisms or skeptical remarks about Semantic Web technology (taken from actual NSF proposal reviews and conference panel remarks).

*“As a constituent technology, ontology work of this sort is defensible. As the basis for programmatic research and implementation, it is a speculative and immature technology of uncertain promise.”*

*“Users will be able to use programs that can understand semantics of the data to help them answer complex questions ... This sort of hyperbole is characteristic of much of the genre of semantic web conjectures, papers, and proposals thus far. It is reminiscent of the AI hype of a decade ago and practical systems based on these ideas are no more in evidence now than they were then.”*

*“Such research is fashionable at the moment, due in part to support from defense agencies, in part because the Web offers the first distributed environment that makes even the dream seem tractable.”*

*“It (proposed research in Semantic Web) pre-supposes the availability of semantic information extracted from the base documents -an unsolved problem of many years, ...”*

*“Google has shown that huge improvements in search technology can be made without understanding semantics. Perhaps after a certain point, semantics are needed for further improvements, but a better argument is needed.”*

**Table 1: Some Reservation among DB researchers about the Semantic Web**

These reservations and skepticism likely stem from a variety of reasons. First, this may be a product of the lofty goals of the Semantic Web as depicted in [Berners-Lee et. al., 2001]. Specifically, database researchers may have reservations stemming from the overwhelming role of description logic in the W3C's Semantic Web Activity and related standards. The vision of the Semantic Web proposed in several articles may seem, to many readers, like a proposed solution to the long standing AI problems. Lastly, one of the major skepticism is related to the legitimate concern about the scalability of the three core capabilities for the Semantic Web to be successful, namely the scalability of (a) ontology creation and maintenance of large ontologies, (b) semantic annotation, and (c) inference mechanisms or other computing approaches involving large, realistic ontologies, metadata, and heterogeneous data sets.

Despite these reservations, some of them well justified, we believe semantic technology is beginning to mature and will play a significant role in the development of future information systems. We believe that database research will greatly benefit by playing critical roles in the development of both Semantic Technology and the Semantic Web. In addition, we also feel that the database community is very well equipped to play their part in realizing this vision. Thus, the aim of this paper is to:

- Identify some prevalent myths about the Semantic Web
- Identify instances of Semantic (Web) Technology in action and how the database community can make invaluable contributions to the same.

For the purpose of this article, as well as for tagging real and existing versus more futuristic and speculative alternatives, we distinguish between Semantic Technology and Semantic Web technology. By Semantic Technology [Polikoff and Allemang 2003] (a term that predates the "Semantic Web"), we imply application of techniques that support and exploit semantics of information (as opposed to syntax and structure/schematic issues [Sheth 99]<sup>1</sup>) to enhance existing information systems. In contrast, the Semantic Web technology (more specifically its vision) is best defined as "*The Semantic Web is an extension of the current web in which information is given well-defined meaning, better enabling computers and people to work in cooperation.*" [Berners-Lee et al 2001]. Currently in more practical terms, Semantic Web technology also implies the use of standards such as RDF/RDFS, and for some OWL. It is however important to note that while description logic is a center piece for many Semantic Web researchers, it is not a necessary component for *many* applications that exploit semantics. For the Semantic Technology as the term is used here, complex query processing, involving both metadata and ontology takes the center piece, and is where the database technology continues to play a critical role. Another term we define for convenience is, semi-formal ontology, based on remarks in [Gruber 2003]. For our purpose, these are ontologies that do not claim formal semantics and/or are populated with partial or incomplete knowledge. For example, in such an ontology, all schema level constraints may not be observed in the knowledgebase that instantiates the ontology schema. This becomes especially relevant when ontology is populated by many persons or by extracting and integrating knowledge from multiple sources.

## **2. Examples of Semantic Technology Applications and Some Observations**

We summarize a few applications developed using commercial technologies to provide insights into what Semantic (Web) Technology can do today. Based on the increasing complexity and the deeper role of semantics, we divide the applications into three types<sup>2</sup>:

- Semantic search and contextual browsing:
  - In *Taalee (now Semagix) Semantic Search Engine* [Townley 2000], the ontology consisted of general interest areas with several major categories (News, Sports, Business, Entertainment, etc.) and over 16 subcategories (Baseball, Basketball, etc in Sports). Blended Semantic Browsing and Querying (BSBQ) provided domain specific search (search based on relevant, domain specific attributes) and contextual browsing. The application involved crawling/extracting audio, video and text content from well over 250 sources (e.g. CNN website). This application was

---

<sup>1</sup> For a commercial use of term "Semantic Technology" see [Polikoff and Allenmang 2003].

<sup>2</sup> At least the applications underlined are known to have been developed by commercial technology/product or deployed.

commercially deployed for a Web-audio company called *Voquette*. An interesting related application not developed by Semagix is reported in [Guha et al 2003].

- Semantic integration:
  - In *Equity Analyst Workbench* [Sheth et al 2002], A/V and text content from tens of sites and NewsML feeds aggregated from 90+ international sources (such as News agencies of various countries) were continuously classified into a small taxonomy, and domain specific metadata was automatically extracted (after one time effort to semi-automatically create a source-specific extractor agent). The equity market ontology used by this application consists of over one million facts (entity and relationship instances). An illustrative example of a complex semantic query involving metadata and ontology this application supported is: Show analyst reports (from many sources in various formats) that are competitors of Intel Corporation.
  - In an application involving *Repertoire Management* for a multinational Entertainment conglomerate, its ontology with relatively simple schema is populated with over 14.5 million instances (e.g., semantically disambiguated names of artists, track names, etc). The application provided integrated access to heterogeneous content in the company's extensive media holding while addressing semantic heterogeneity.
- Analytics and Knowledge Discovery:
  - In the *Passenger Threat Assessment* application for national/homeland security [Sheth et al 2004] and Semagix's *Anti-money Laundering* solution [Semagix-CIRAS], the knowledge base is populated from many public, licensed and proprietary knowledge sources. The resulting knowledge base has over one million instances. Periodic or continuous metadata extraction from tens of heterogeneous sources (150 files formats, HTML, XML feeds, dynamic Web sites, relational databases, etc) is also performed. When the appropriate computing infrastructure is used, the system is scalable to hundreds of sources, or about a million documents per day per server. A somewhat related non-Semagix business intelligence [IBMWF] application has demonstrated scalability by extracting metadata (albeit somewhat limited types of metadata with a significantly smaller ontology) from a billion pages [Dill et al 2003].

Based on our experience in building the above real-world applications, we now review some empirical observations:

1. Applications validate the importance of ontology in the current semantic approaches. An ontology represent a part of the domain or the real-world for which it represents and captures a shared knowledge around which the semantic application revolves. It is the "ontological commitment" reflecting agreement among the experts defining the ontology and its uses, that is the basis for "semantic normalization" necessary for semantic integration.
2. Ontology population is critical. Among the ontologies developed by Semagix or using its technology, median size of ontology is over 1 million facts. This level of capture of knowledge makes the system very powerful. Since it is obvious that this is the sort of scale Semantic Web applications are going to be dealing with, means of populating ontologies with instance data need to be automated.
3. Two of the most fundamental "semantic" techniques are named entity, and semantic ambiguity resolution (Also closely tied to data quality problem). Any semantic technology and its application. Without good solutions to these none of the applications listed will be of any practical use. For example, a tool for annotation is of little value if it does not support ambiguity resolution. Both require highly multidisciplinary approaches, borrowing for NLP/lexical analysis, statistical and IR techniques and possibly machine learning techniques.
4. Semi-formal ontologies that may be based on limited expressive power are most practical and useful. Formal or semi-formal ontologies represented in very expressive languages (compared to moderately expressive ones) have in practice, yielded little value in real-world applications. One reasons for this may be that it is often very difficult to capture the knowledge that uses the more expressive constructs of a representation language. This difficulty is especially apparent when trying to populate an ontology that uses a very expressive language to model a domain. Hence the additional effort in modeling these constructs for a particular domain is often not justifiable in terms of the gain in performance. Also there is widely accepted trade-off between expressive power and computational complexity associated with inference mechanisms for such languages. Practical applications often end up using languages that lie closer to less expressive languages in the "expressiveness vs. computational complexity continuum". This resonates with so-called Hendler's hypothesis ("little semantics goes a long way").

5. Large scale metadata extraction and semantic annotation is possible. Storage and manipulation of metadata for millions to hundreds of millions of content items requires best applications of known database techniques with challenge of improving upon them for performance and scale in presence of more complex structures.
6. Support for heterogeneous data is key – it is too hard to deploy separate products within a single enterprise to deal with structured and unstructured data/content management. New applications involve extensive types of heterogeneity in format, media and access/delivery mechanisms (e.g., news feed in NewsML news, Web posted article in HTML or served up dynamically through database query and XSLT transformation, analyst report in PDF or WORD, subscription service with API-based access to Lexis/Nexis, etc). Database researchers have long studied the issue of integrating heterogeneous data, and many of these come handy.
7. Semantic query processing with the ability to query both ontology and metadata to retrieve heterogeneous content is highly valuable. Consider the query “Give me all articles on the competitors of Intel”, where ontology gives information on competitors, supports semantics (with the understanding that “Palm” is a company and that “Palm” and “Palm, Inc.” are the same in this case), and metadata identifies the company an article refers to, regardless of format of the article. Analytical applications could require sub-second response time for tens of concurrent complex queries over large metadata base and ontology, and can benefit from further database research. High performance and highly scalable query processing that deal with more complex representations compared to database schemas and with more explicit role of relationships, is important. Database researcher can also contribute to the strategies of dealing with large RDF stores.
8. A vast majority of the Semantic (Web) applications that have been developed or envisioned rely on three crucial capabilities namely ontology creation, semantic annotation and querying/inferencing. Enterprise scale application share many requirements in these three respects with pan Web applications. All these capabilities must scale to millions of documents and concepts (rather than hundreds to thousands). Main differences are in the number of content sources and the corresponding size of metadata.

### 3. Discussion

Ontologies come in bewildering variety; Figure 1 represents just three of the dimensions. To keep a focus on real world applications and for the sake of brevity, we restrict the scope to task specific and domain specific ontologies. As observed recently by Gruber [Gruber 2003], currently the ontologies that are semi-formal have demonstrated very high practical value. We believe ontology development effort for semi-formal ontologies can be significantly smaller, especially for the ontology population effort, compared to that required for developing formal ontologies or ontologies with more expressive representations. Semi-formal ontologies have provided good examples of both value and utility.

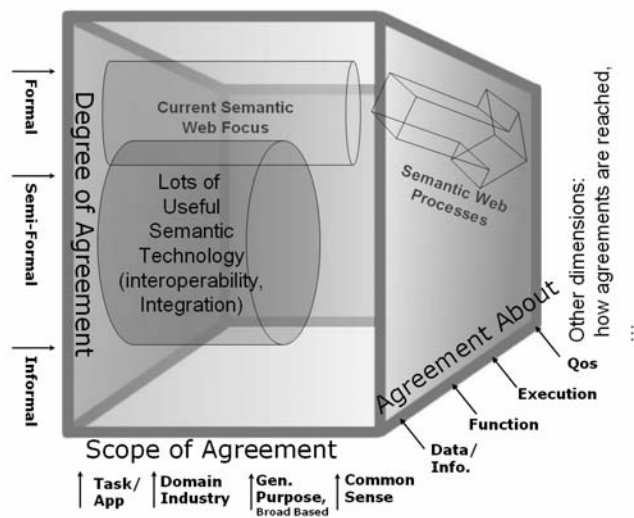


Figure 1 Dimension along which ontologies vary

For example, GO (<http://www.geneontology.org/>), is arguably a nomenclature from the perspective of representation and lacks formal and richer representation necessary to qualify as a formal ontology (this is discussed in more detail later). Research in database and information systems have played and will continue to play a critical role with respect to the scalability. We review the crucial scalability aspect of the three capabilities next.

## Availability of large and “useful” ontologies

Although an ontology schema may resemble at a representational level a database schema, and instances may reflect database tuples, the fundamental difference is that ontology is supposed to capture some aspect of real-world or domain semantics, as well as represent ontological commitment forming the basis of semantic normalization. Methods for creating ontologies can be grouped into the following types:

- Social processes where a group of users go through a process of suggestions and iteratively revise versions of ontologies to capture domain semantics.
- Automatically (or semi-automatically) extract the ontology schema (i.e., ontology learning) from content of various kinds. Although there has been a recent spate of interest, this approach relates to the knowledge acquisition bottleneck faced in the AI research of eighties, and we have little practical experience to rely upon [Gómez-Pérez and Manzano-Macho 2003].
- Automatic or semi-automatic (with human curation) population of the knowledge base with respect to human design ontology schema. We can report on practical experience with a scalable approach of this type since several ontologies with over million instances have been create with a total of 4 to 8 weeks of effort (e.g., knowledge extraction in *SemagixFreedom* [Sheth et. al., 2002]).

## Semantic Metadata Extraction or Semantic Annotation of massive content

Annotating heterogeneous content with semantics provided by relevant ontology (or ontologies) has been identified as a key challenge for the Semantic Web. Recently there have been commercial results providing detailed semantic annotations of heterogeneous content (structured, semi-structured, and unstructured with different formats) [Sheth et al 2002, Hammond et al 2002], as well as research reporting annotation of over a billion Web pages [Dill et al 2003]. As observed in the efforts on automatic semantic annotation, two resources necessary for realizing the semantic web are: (a) large scale availability of domain specific ontologies; and (b) scalable techniques to annotate content with high quality metadata descriptions based on the terms, concepts or relationships provided by these ontologies. We believe main area of challenge here is to support increasing number of practical and scalable techniques for semantic disambiguation.

## Inference Mechanisms that Scale

Inference mechanisms that can deal with the massive number of assertions that would be encountered by Semantic Web applications are required. The claimed power behind many of the proposed applications of Semantic Web technology is the ability to infer knowledge that is not explicitly expressed. Needless to say, this feature has attracted attention from the AI community since they have been dealing with issues relating to inference mechanisms for quite some time. Inference mechanisms are applicable only in the context of formal ontologies. They use rules and facts to assert new facts that were not previously stated as true. One of the most common knowledge representation languages has been Description Logic [Nardi et. al., 2002] on which DAML, one of the earliest Semantic Web languages is based. It was in fact one of the less expressive members of the DL family. The reason for limiting the expressive power of such a knowledge representation formalism was very clear when the decidability and complexity issues were considered. Although several optimized methods of inference were introduced later [Baader et.al., 2001], inference mechanisms were still overshadowed by the performance advantage of traditional database systems. This has lead to reluctance among many database researchers to accept the Semantic Web vision as viable. Description Logics may form a part of some Semantic Web solutions of the future. We are however convinced it is not the only knowledge representation formalism that will go on and make the Semantic Web vision a reality. One may argue that it is possible to do some sort of rudimentary inference using RDFS (using *subClassOf* and *subPropertyOf*). However, using RDF/RDFS does not force one to use *only* inference mechanisms of some sort in applications. Since RDFS has a graph model associated with it there is the possibility to use other techniques to answer complex queries [Sheth et. al., 2004].

## Why semi-formal, less expressive ontologies?

Ontologies serve several purposes, including: having an agreement between humans, having a common representation for knowledge, having machines (software) get common interpretation of something that humans have agreed to, and forming the basis for defining metadata or semantic annotation. Tom Gruber, who many would credit with bringing the term to vogue in contemporary knowledge representation research, identifies three types of ontologies—informal, semiformal and formal [Gruber 2003]. He stresses

the value of semi-formal ontology in meeting several challenges; especially that of information integration. Some researchers in the Semantic Web community would argue for only formal ontologies (and discount the value of semi-formal ontologies). We do not doubt that formal ontologies have a potential role in Semantic Web research. However, database researchers should particularly realize the value of and exploit semi-formal ontologies. **Figure 1** stresses that there is a very large body of work that can and needs to be done using semi-formal ontologies.

There is a very good reason as to why semi-formal ontologies are both more abundant and more useful than formal ontologies. The answer lies in the ease with which semi-formal ontologies can be built to a scale that is useful in real-world applications. One key reason is that of the need to accommodate partial (incomplete) and possibly inconsistent information, especially in the assertions of an ontology. This view is consistent with the view presented in [Shirky 2003] as (replacing “standard” with “ontology”): “the more semantic consistency required by a standard, the sharper the tradeoff between complexity and scale.” GO ontology, which is more a nomenclature and taxonomy, than a formal ontology, is highly successful and extensively used. Although GO is technically a nomenclature rather than an ontology [Kremer2002]<sup>3</sup>, it has been shown to have several inconsistencies, it has been successfully used to annotate large volumes of data and consequently support interoperability and integration from heterogeneous data sets. This shows that highly expressive formal ontologies are not required for *all* Semantic Web applications. It also shows that real world applications often can be developed with very little semantics (cf: Jim Hendler’s hypothesis: “little semantics goes a long way”), or with compromises with completeness and consistency required by more formal representations and inferencing techniques.

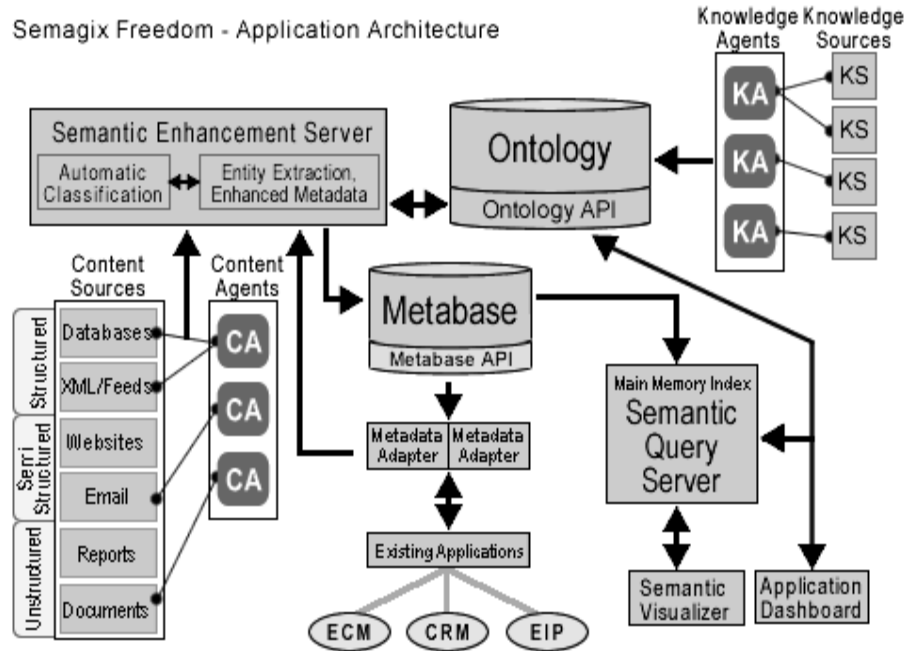
Our objective in touting the value of semi-formal ontologies is to prevent research in the Semantic web field from leading straight into the very problems that AI found itself in. We hope to do this by reducing the prevalent emphasis on formal ontologies and pure deductive inference mechanisms. The reader should note however that we do not completely discount the value of the same. We liken some of the current research direction in the Semantic Web community to, attempting to construct a new building using the flawed building blocks that lead to the downfall of previous building attempts. Our reasoning behind this is that most motivating examples described in this field pay little or no attention to the fundamental (read hard) problems of entity/relationship identification and ambiguity resolution. Database researchers working on schema integration are only too familiar with the problems relating to ambiguity resolution. According to [Shirky 2003], most scenarios described for potential applications of the Semantic Web trivialize these fundamentally hard problems while emphasizing the trivial problems. Our views coincide with those expressed in [Kremer 2002, Brodie 2003]. In [Brodie 2003] the Semantic Web community is urged to not waste their efforts on “fixing the plumbing” (referring to infrastructure issues) and to focus their efforts on the more fundamental problems.

#### **4. Semagix Freedom: An example of state of the art Semantic Technology**

Let us briefly describe a state of the art commercial technology and product that is built upon the key perspectives we presented above. Semagix Freedom exploits task and domain ontologies that are populated with relevant facts in all key functions: automatic classification of content, ontology-driven metadata extraction, and support for complex query processing involving metadata and ontology for all three types of semantic applications identified in Section 2. It provides tools that enable automation in every step in the content chain - specifically ontology design, content aggregation, knowledge aggregation and creation, metadata extraction, content tagging and querying of content and knowledge. Scalability, supported by a high degree of automation and high performance based on main memory based query processing has been of critical importance in building this commercial technology and product. Figure 2 below shows the architecture of Semagix Freedom.

---

<sup>3</sup> “For data annotation, in principle not a full fledged ontology as described above is required but only a controlled vocabulary since the main purpose is to provide constant and unique reference points.”[Kremer2002]

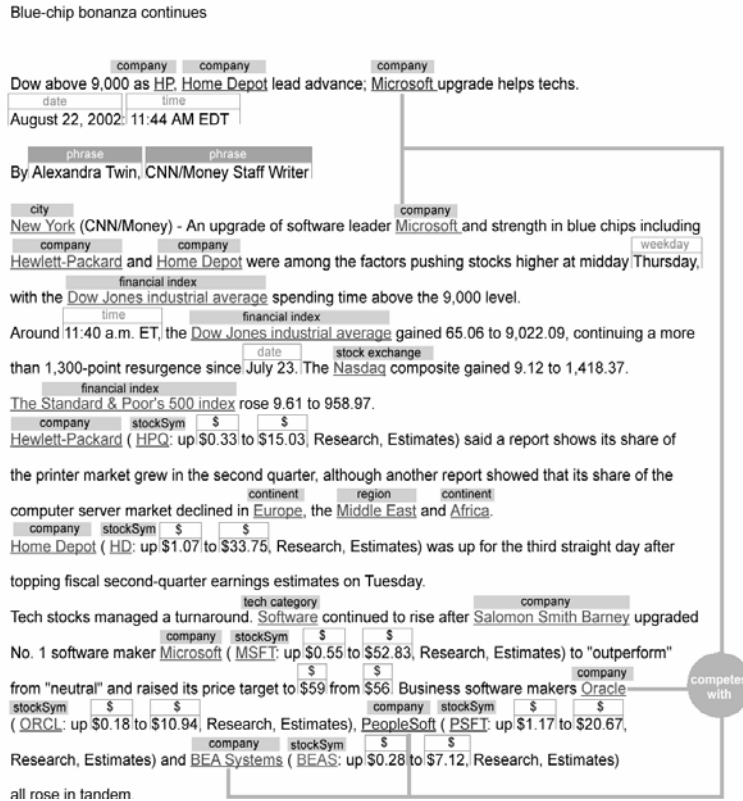


**Fig. 2. Semagix Freedom Architecture**

Freedom provides a modeling tool to design the ontology schema based on the application requirements. The domain specific information architecture is dynamically updated to reflect changes in the environment, and it is easy to configure and maintain. The Freedom ontology is populated with knowledge, which is any factual, real-world information about a domain in the form of entities, relationships, attributes and certain constraints. The ontology is automatically maintained by Knowledge Agents (Figure 2, top right). These are software agents created without programming that traverse trusted knowledge sources that may be heterogeneous, but either semi-structured or structured (i.e., concept extraction from plain text to populate ontology is currently not supported but may be supported in future). Knowledge Agents exploit structure to extract useful entities and relationships for populating the ontology automatically. Once created, they can be scheduled to automatically keep the ontology up-to-date with respect to changes in the knowledge sources. Semantic ambiguity resolution (is the entity instance the same or related to an existing entity instance? Is this the same “John Doe” Board Member the same as the “John Doe” CEO in the ontology) is one of the most important capabilities associated with this activity, as well as with the metadata extraction. Ontology can be exported in RDF/RDFS barring some constraints that cannot be presented in RDF/RDFS.

Freedom also aggregates structured, semi-structured and unstructured content from any source and format. Two forms of content processing are supported: automatic classification and automatic metadata extraction. Automatic classification utilizes a classifier committee based on statistical, learning, and knowledgebase classifiers. Metadata extraction involves named entity identification and semantic disambiguation to extract syntactic and contextually relevant semantic metadata (Figure 2, left). Custom meta-tags, driven by business requirements, can be defined at a schema level. Much like Knowledge Agents, Content Agents are software agents created without programming using an extensive toolkit. Incoming content is further “enhanced” by passing it through the Semantic Enhancement Server [Hammond et al 2002]. The Semantic Enhancement Server can identify relevant document features such as currencies, dates, etc., perform entity disambiguation, tag the metadata with relevant knowledge (i.e., the instances within the ontology) and produce a semantically annotated content (that references relevant nodes in the ontology) or a tagged output of metadata. Automatic classification aid metadata extraction and enhancement by providing context needed to apply the relevant portion of a large ontology.





**Figure 3: An example of Automatic Semantic Metadata Extraction/Annotation**

The Metabase stores both semantic and syntactic metadata related to content. It stores content into a relational database as well as a main-memory checkpoint. At any point in time, a snapshot of the Metabase (index) resides in main memory (RAM), so that retrieval of assets is accelerated using the Semantic Query Server. This index is both incremental (to keep up with new metadata acquisition) and distributed (i.e., layered over multiple processors, to scale with number of contents and size of the Metabase). The Semantic Query Server is a main memory–based front–end query server. The Semantic Enhancement and Query Servers provide semantic applications (or agents) ability to query Metabase and ontology using http and Java-based APIs, returning, returning results in XML with published DTDs. This ability, with the context provided by ontology and ambiguity resolution, form the basis for contextual, complex, and high performance query processing, providing highly relevant content to the semantic applications . Let's end the review of Freedom by summarizing some of its scalability and performance capabilities *to date*, along with some experiences based on development of Semantic Applications for paying (i.e., real-world) customers:

1. Typical size of an ontology schema for a domain or task ontology: 10s of (entity) classes, 10s of relationships, few hundred property types
2. Average size of ontology population (number of instances): over a million of named entities
3. Number of instances that can be extracted and stored in a day (before human curation, if needed): up to a million per server per day
4. Number of text documents that can be processed for automatic metadata extraction per server per day: hundreds of thousand to a million
5. Performance for search engine type keyword queries: well over 10 million queries per hour with approx. 10ms per query for 64 concurrent users
6. Query processing requirement observed in an analytical application: approx. 20 complex queries (involving both Ontology and Metabase) to display a page with analysis, taking a total of 1/3 second for computation (roughly equivalent to 50+ query over a relational database with response time over 50 seconds).

## 5. Conclusion

Formal ontologies in description logic based representation; supported by deductive inference mechanisms may not be the primary (and certainly not the only) means of addressing major challenges in realizing the Semantic Web vision. The database community should realize that the Semantic Web vision is not one of solving the AI problem, or OWL with subsumption based inference mechanisms. Instead, it can make critical contribution to the Semantic Web by drawing upon its past work and further research on topics such as supporting processing of heterogeneous data/content, semantic ambiguity resolution, complex query processing involving metadata and knowledge represented in semi-formal ontology, and ability to scale with large amount of structured and semi-structured information. In supporting this view point, we provided an overview of one instance of the commercial technology that has been used to develop a broad variety of real world semantic applications. We also provided high level information on scalability requirements observed in supporting these applications. Alternative strategies to realize the vision of the Semantic Web will need to show they will need to scale and perform at least as well as what today's commercial technologies (such as the one briefly discussed in this article) do, and probably well beyond that.

## Acknowledgements

Comments to a draft from Dean Allemang, Kemafor Anyanwu, Vipul Kashyap, Harumi Kuno and Kevin Wilkinson are gratefully acknowledged. We also received editing assistance and comments from members of the LSDIS lab including Chris Halaschek, Christopher Thomas, Meenakshi Nagarajan and Kunal Verma.

## References

- [Baader et. al., 2001] F. Baader and U. Sattler. An Overview of Tableau Algorithms for Description Logics. *Studia Logica*, 69:5-40, 2001.
- [Berners-Lee et. al., 2001] T. Berners-Lee, J. Hendler and O. Lassila. The Semantic Web A new form of Web content that is meaningful to computers will unleash a revolution of new possibilities, *Scientific American*, May 2001.
- [Brodie 2003] M. Brodie. The Long and Winding Road To Industrial Strength Semantic Web Services, Keynote Talk, ISWC 2003. <http://iswc2003.semanticweb.org/brodie.pdf>
- [Dill et. al., 2003] S. Dill et. al. SemTag and SemSeeker: Bootstrapping the Semantic Web via automated semantic annotation. Proceedings of the 12th International WWW Conference (WWW 2003), Budapest, Hungary, May 2003.
- [Fisher and Sheth 2003] M. Fisher and A. Sheth. "Semantic Enterprise Content Management," Practical Handbook of Internet Computing, Munindar P. Singh, Ed., CRC Press, 2003.
- [Gómez-Pérez and Manzano-Macho 2003] A. Gómez-Pérez, D. Manzano-Macho, A survey of ontology learning methods and techniques, 2003.
- [Gruber 2003] T. Gruber. It Is What It Does: The Pragmatics of Ontology, invited talk at Sharing the Knowledge- International CIDOC CRM Symposium, March 26-27, Washington, DC, <http://tomgruber.org/writing/cidoc-ontology.htm>
- [Guha et. al., 2003] R. Guha, Rob McCool and Eric Miller. Semantic Search, The Twelfth International World Wide Web Conference, Budapest Hungary, May 2003
- [Hammond et. al., 2002] B. Hammond, A. Sheth, and K. Kochut. Semantic Enhancement Engine: A Modular Document Enhancement Platform for Semantic Applications over Heterogeneous Content, in *Real World Semantic Web Applications*, V. Kashyap and L. Shklar, Eds., IOS Press, December 2002, pp. 29-49
- [Handschuh et. Al., 2003] S. Handschuh, S. Staab and R. Volz. On Deep Annotation. Proceedings of the 12th International WWW Conference (WWW 2003), Budapest, Hungary. May 2003.
- [IBM-WF] WebFountain, [http://www-1.ibm.com/mediumbusiness/venture\\_development/emerging/wf.html](http://www-1.ibm.com/mediumbusiness/venture_development/emerging/wf.html)
- [Kahn et. al., 2001] J. Kahan, M-R. Koivunen, E. Prud'Hommeaux and R. Swick. Annotea: An open RDF Infrastructure for shared annotations. Proceedings of the 10th International WWW Conference (WWW 2002), Hong Kong, May 2001.

Slightly abridged version appears in *IEEE Data Engineering Bulletin, Special issue on Making the Semantic Web Real*, U. Dayal, H. Kuno, and K. Wilkinson, Eds. December 2003.

- [Kremer 2002] S. Schulze-Kremer: Ontologies for molecular biology and bioinformatics. In *Silico Biology* 2: 17, 2002.
- [McDermott 1987] D. McDermott, Critique of Pure Reason Computational Intelligence, 3:151-237, 1987.
- [Polikoff and Allemang 2003] I. Polikoff and D. Allemang, "Semantic Technology," TopQuadrant Technology Briefing v1.1, September 2003.  
[http://www.topquadrant.com/documents/TQ03\\_Semantic\\_Technology\\_Briefing.PDF](http://www.topquadrant.com/documents/TQ03_Semantic_Technology_Briefing.PDF)
- [Semagix-CIRAS] Anti-Money Laundering, Semagix, Inc. [http://www.semagix.com/solutions\\_ciras.html](http://www.semagix.com/solutions_ciras.html)
- [Sheth 1999] A. Sheth, "Changing Focus on Interoperability in Information Systems: From System, Syntax, Structure to Semantics", in *Interoperating Geographic Information Systems*. M. F. Goodchild, M. J. Egenhofer, R. Fegeas, and C. A. Kottman (eds.), Kluwer, Academic Publishers, 1999, pp. 5-30.
- [Sheth and Meersman 2002] A. Sheth and R. Meersman, "Amicalola Report: Database and Information Systems Research Challenges and Opportunities in Semantic Web and Enterprises," *ACM SIGMOD Record*, Vol. 31, No. 4, December 2002, pp. 98-106. <http://lstdis.cs.uga.edu/SemNSF/>
- [Sheth et. al., 2002] A. Sheth, C. Bertram, D. Avant, B. Hammond, K. Kochut and Y. Warke. Semantic Content Management for Enterprises and the Web, *IEEE Internet Computing*, July/August 2002, pp. 80-87.
- [Sheth et. al., 2004] A. Sheth, et al. Semantic Association Identification and Knowledge Discovery for National Security Applications, *Journal of Database Management*, 2004 (to appear).
- [Shirky 2003] C. Shirky. The Semantic Web, Syllogism, and Worldview. In *Networks, Economics, and Culture*, November 7, 2003.
- [Townley 2000] J. Townley, The Streaming Search Engine That Reads Your Mind, *Streaming Media World*, August 10, 2000. <http://smw.internet.com/gen/reviews/searchassociation/index.html>.
- [Nardi et. al., 2002] D. Nardi, R. J. Brachman. An Introduction to Description Logics. In the *Description Logic Handbook*, edited by F. Baader, D. Calvanese, D.L. McGuinness, D. Nardi, P.F. Patel-Schneider, Cambridge University Press, 2002, pages 5-44.