

Next generation knowledge management

J Davies, R Studer, Y Sure and P W Warren

Despite its explosive growth over the last decade, the Web remains essentially a tool to allow humans to access information. The next generation of the Web, dubbed the 'Semantic Web', will extend the Web's capability through the increased availability machine-processable information. These machine-processable descriptions of Web information resources are called meta-data and are associated with ontologies, or conceptualisations of the domain of application. Meta-data and associated ontologies then allows more intelligent software systems to be written, automating the analysis and exploitation of Web-based information.

This paper describes how knowledge management can be improved through the adoption of Semantic Web technology. To realise this, a number of different technologies need to be brought together. Their fusion provides the infrastructure which makes semantic knowledge management possible. Specifically, the paper discusses the use of knowledge discovery and human language technology to (semi-)automatically derive the required ontologies and meta-data, along with a methodology to support this process. We describe techniques for management and controlled evolution of ontologies and a set of semantic knowledge access tools for enhanced information access. Finally, a set of application scenarios for the technology are sketched.

1. Introduction

There are now many tens of billions of documents on the WWW, which are used by more than 300 million users globally, and millions more pages on corporate intranets and extranets. The continued rapid growth in information volume makes it increasingly difficult to find, organise, access and maintain the information required by users. Contemporaneously with this explosion of Web-based information, the notion of a semantic Web [1] has been proposed that has the potential to provide enhanced information access based on the exploitation of machine-processable meta-data. In this paper, we are particularly interested in the new possibilities afforded by Semantic Web technology in the area of knowledge management.

Until comparatively recently, the value of a company was felt to be determined mainly by the value of its tangible assets. In recent years, however, it has been increasingly recognised that in the post-industrial era, an organisation's success is more dependent on its intellectual assets than on the value of its physical resources.

The requirement for highly skilled labour in many industries, new computing and telecommunications technologies, faster innovation, and ever shorter product cycles, has caused a huge change in the ways

organisations compete — knowledge is now the key battleground for competition.

Other factors driving companies to try and manage and exploit their intellectual assets more effectively are:

- increasing employee turnover rates and a more mobile workforce, which can lead to loss of knowledge,
- globalisation, often requiring people to collaborate and exchange knowledge across continents and time zones.

The knowledge management (KM) discipline aims to address this challenge and can be broadly defined as the tools, techniques and processes for the most effective and efficient management of an organisation's intellectual assets [2]. These intellectual assets can be exploited in a variety of ways. By sharing and reusing current best practice, for instance, you can improve current business processes and eliminate duplication of effort. New business opportunities can be generated by collecting intelligence on markets and sales leads; and new products and services can be created, developed and brought to the market-place ahead of competitors.

It has often been argued in KM circles that technology is a relatively marginal aspect of any KM

initiative and that organisational culture is a more important feature. While the sentiment that we need a wider perspective than just technology is correct, this argument reveals the assumption of a dichotomy between technology and organisational culture which does not exist. Rather, technology-based tools are among the many artefacts entwined with culture, whose use both affects and is affected by the prevailing cultural environment. A holistic view is required and technology often plays a larger part in cultural factors than is sometimes acknowledged¹.

The focus of this paper is research into semantic Web-based tools for knowledge management; it is, however, equally important to understand the cultural and organisational contexts in which such tools can be used to best effect. Related work in this area can be found in, for example, Antoniou and van Harmelen [3].

2. The Semantic Web and knowledge management

Intranets have an important role to play in the more effective exploitation of both explicit (codified) and tacit (unarticulated) knowledge. With regard to explicit knowledge, intranet technology provides a ubiquitous interface to an organisation's knowledge at relatively low cost using open standards. Moving information from paper to the intranet can also have benefits in terms of speed of update and hence accuracy. The issue then becomes how to get the right information to the right people at the right time — indeed, one way of thinking about explicit knowledge is that it is information in the right context, i.e. information which can lead to effective action. Regarding tacit knowledge, technology also has a role to play, since we can use intranet-based tools to connect people with similar interests or concerns, thus encouraging dialogue and opening up the possibility of the exchange of tacit knowledge.

Important information is often scattered across Web and/or intranet resources. Traditional search engines return ranked retrieval lists that offer little or no information on the semantic relationships between documents. Knowledge workers consequently spend a substantial amount of their time browsing and reading to find out how documents are related to one another and where each falls into the overall structure of the problem domain. Yet only when knowledge workers begin to locate the similarities and differences between pieces of information do they move into an essential part of their work — building relationships to create new knowledge.

¹To take an obvious example, consider the way in which the widespread introduction of e-mail over the last decade or so has changed ways of working.

Current knowledge management systems have significant weaknesses.

- Searching information

Existing keyword-based searches can retrieve irrelevant information that includes certain terms in different meanings. They also miss information when different terms with the same meaning about the desired content are used. Information retrieval traditionally focuses on the relationship between a given query (or user profile) and the information store. On the other hand, exploitation of interrelationships between selected pieces of information (which can be facilitated by the use of ontologies) can put otherwise isolated information into a meaningful context. The implicit structures so revealed help users use and manage information more efficiently.

- Extracting information

Currently, human browsing and reading is required to extract relevant information from information sources. This is because automatic agents do not possess the commonsense knowledge required to extract such information from textual representations, and they fail to integrate information distributed over different sources.

- Maintenance

Maintaining weakly structured text sources is a difficult and time-consuming activity when such sources become large. Keeping such collections consistent, correct, and up-to-date requires mechanised representations of semantics that help to detect anomalies.

- Automatic document generation

This would enable adaptive Web sites that are dynamically reconfigured according to user profiles or other aspects of relevance. Generation of semi-structured information presentations from semi-structured data requires a machine-accessible representation of the semantics of these information sources.

To sum up, we want to move from a document-centric view of information retrieval to a knowledge-centric view, wherein tools are not returning ranked lists of documents to the user, but, instead, attempt to provide them with the specific information they need perhaps gathered from multiple documents.

Knowledge management tools are needed that can integrate the resources dispersed across Web resources into a coherent corpus of interrelated information. Previous research in information integration has largely focused on integrating heterogeneous databases and

knowledge bases, which represent information in a highly structured way, often by means of formal languages. In contrast, the Web consists to a large extent of unstructured or semi-structured natural language text.

The goal of the Semantic Web is to offer automated information access based on machine-processable semantics of data and heuristics that use these semantics.

The explicit representation of the semantics of data, accompanied with domain theories (i.e. ontologies — see Fig 1), will enable a Web that provides a qualitatively new level of service. It will weave together a large network of human knowledge and will complement it with machine processability. Various automated services will help the user achieve goals by accessing and providing information in machine-understandable form.

Ontologies offer a way to cope with heterogeneous representations of Web resources. The domain model implicit in an ontology can be taken as a unifying structure for giving information a common representation and semantics. A semantic repository of information spanning multiple heterogeneous information sources can be created.

The use of ontologies and supporting tools offer an opportunity to significantly improve knowledge management capabilities in large organisations, and it is their use in this particular area which is the subject of this paper.

3. Creating the semantic repository — ontologies and instances

3.1 *Ontology learning*

To date, much ontology creation has been a manual process. In the CYC [4] project, for example, common-sense knowledge was extracted manually from different sources and expressed using ontologies. A similar approach is used by Yahoo and also by Google Directory, which is based on the dmoz open directory project [5]. This is inevitably a very labour-intensive process, and there is a need to at least partially automate it. The need is to identify some classes and instances automatically. We can imagine a predefined ontology of classes and relationships, plus a knowledge base of instances, being extended by automated learning. Alternatively, an ontology and knowledge base might be learned from scratch, albeit with the need for some human intervention. Most work so far has been in the former category.

The most promising approach for ontology learning applies knowledge discovery techniques, based on statistics and machine learning, to text mining. Indeed the machine learning community is used to employing models in their automated learning algorithms. Ontologies are simply another class of model, although rather more complex than normally used in machine learning. A typical approach uses clustering and related techniques. For example, work has been done to extend an existing WordNet ontology [6].

A standard approach to performing document, or web page, clustering is to regard each document as a

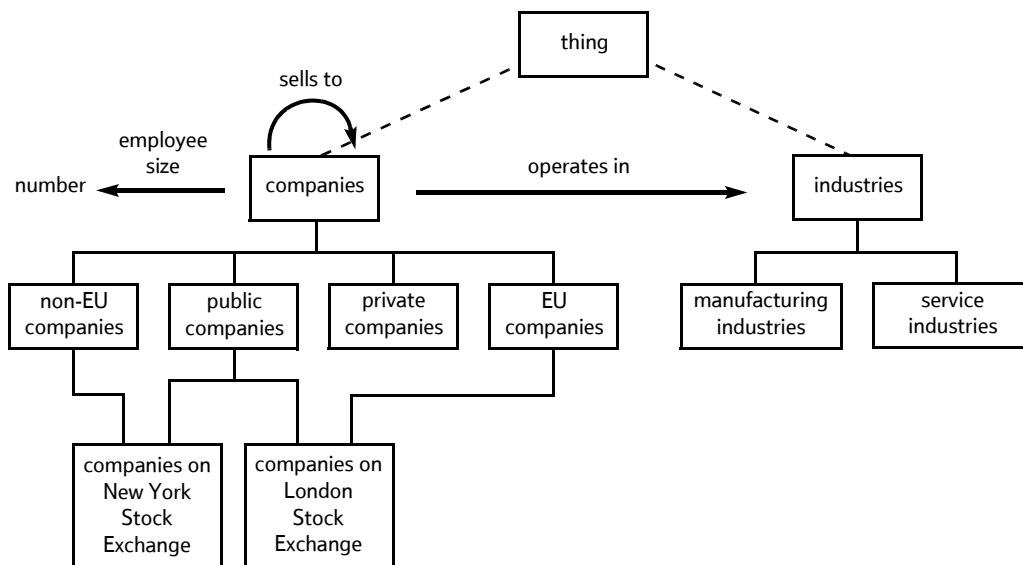


Fig 1 An example of a part of an ontology.

vector in a high dimensional space [7]. The dimensionality of the vector space is defined by the words used in the whole corpus of documents. Any document is thereby described by a vector listing the number of occurrences of each word in the document. This enables the computation of a similarity measure between two documents X and Y , with components X and Y respectively. One of the most obvious is cosine similarity², calculated by:

$$\cos(X, Y) = \frac{\sum_i X_i Y_i}{\sqrt{\sum_j X_j^2 \sum_l Y_l^2}}$$

In vector algebra, this is the normalised dot product of the two vectors. In the extreme, and highly unlikely situation, where the two documents used exactly the same words with exactly the same frequency, this measure would be unity. At the other extreme where there were no words in common, the measure would be zero. Given such a similarity measure, it is possible to identify which documents cluster closely together [8].

Using cluster analysis to identify classes of documents is perhaps the most obvious form of ontology learning. It is easy to see that, given these classes (clusters), and given a similarity measure, it is possible to associate new instances (documents) with one or more classes. Moreover, clustering can be used to create a hierarchy of classes. In this case it may be appropriate to associate a document with an ultimate class, i.e. leaf node, or with an intermediate node. Document categorisation algorithms are able to accommodate these possibilities.

The same techniques can be used to learn classes within documents. For example, we can work at the level of sentences. Each sentence within a document is represented as a word vector, sentences are clustered, and each cluster is labelled by the most characteristic words from its sentences. Some researchers have used WordNet to improve the results by mapping these clusters to the classes of the general WordNet ontology [9]. These found classes are then used as semantic labels, i.e. XML tags, for annotating documents.

In a given application, the total automation of ontology learning may not yield sufficiently good results. We expect that human intervention and guidance will be required, e.g. in the form of selecting or rejecting suggested classes.

² Various other, more sophisticated approaches have been proposed in the literature including, for example, latent semantic indexing.

We have already observed that, so far, much ontology creation has been a manual process. In this way extensive ontologies have been created for particular specialities, e.g. for the Medline collection of medical papers. This suggests that the human expertise captured in these existing ontologies can be extracted using knowledge discovery techniques. This has, for example, been applied to the Yahoo topic ontology [10].

3.2 Entity and relationship extraction

The techniques discussed in the last section can be used to identify classes and associate documents or parts of documents with those classes. Another technology, that of information extraction, can be used to identify named entities within text, to associate those named entities with classes, and to identify relationships between the entities. The most obvious, and most common, form of a named entity is a proper noun, written in English with a capital letter. This might be the name of a person ('George W Bush'), of a geographical entity ('USA') or of a role ('President'). Such named entities can be regarded as instances of an ontology. With machine learning it is possible, with a reasonable degree of accuracy, to relate the instances to their classes, i.e. person, country, role in the above example.

With information extraction technology it is also possible to perform what is called 'co-reference resolution', i.e. where different named entities refer to the same underlying instance. In one article we might have 'George W Bush', 'Mr. Bush', 'the President', 'he', 'him' all referring to the same person. There are now algorithms capable of good results in co-reference resolution. Some equivalences are easier to resolve than others. Establishing identity between names with different spellings (e.g. 'BT', 'British Telecom', 'British Telecommunications plc') is significantly easier than establishing identity between pronouns and names. Words such as pronouns which refer back to other words or phrases are termed anaphora, and establishing the equivalence of anaphora to other linguistic units is termed anaphora resolution.

In the knowledge management area, it is easy to see the implications that co-reference resolution has for search technology, for example. Many search engines today, when requested to search for a particular string, do not just return the appropriate documents, but also highlight the occurrences of the string within the document. Using co-reference resolution, a search for 'George W Bush' would also highlight locations in the documents where 'the President', 'Mr Bush', occurred, or might even highlight 'he' or 'him' where the pronoun referred to George Bush.

The field of information extraction existed prior to the beginnings of the semantic Web. An important

driver for its development has been the demands of the intelligence community in the USA. Like knowledge discovery, it is now being directed towards ontologies. Named entity extraction and co-reference resolution can be used to identify instances and establish where differing terms refer to the same instance.

Besides entity extraction and coreference resolution, other functions supported by information extraction, and which have relevance for identifying and learning about instances, are:

- associating descriptions with named entities, examples of this perhaps being identifying that 'Bush administration' is of type 'government' and 'Washington' is of type 'city' — in semantic knowledge technology this can be used to associate entities with classes,
- learning relations between the entities — to-date, these have been typically relations such as 'located in' and 'works for',
- identifying events, each of which is essentially a set of entities and relationships.

Figure 2 illustrates the five basic operations of information extraction, applied to a simple text.

Ryanair announced yesterday that it will make Shannon its next European base expanding its route network to 14 in an investment worth around €180m. The airline says it will deliver 1.3 million passengers in the first year of the agreement, rising to two million by the fifth year.

- **entities:** Ryanair, Shannon
- **mentions:** it = Ryanair, the airline = Ryanair, it = the airline
- **descriptions:** European base
- **relations:** Shannon base_of Ryanair
- **events:** investment (€180m)

Fig 2 Extracting information.
(Courtesy: University of Sheffield [11])

3.3 Creating a knowledge base

Along with the ontology, containing the classes and properties, we also have a 'knowledge base' which contains the instances and the specific property instantiations³. Together the ontology and the knowledge base make up a semantic repository.

³Readers should be aware that some authors do not maintain the sharp distinction made here between ontology and knowledge base, and use the term 'ontology' to include the instances and property instantiations.

Whether the two parts are stored separately, e.g. in two distinct relational databases, depends on the practicalities of the implementation.

As a text is analysed and named entities identified, hyperlinks are established to the instances in the knowledge base. Figure 3 illustrates a piece of text and an associated semantic repository. The repository contains both ontology (classes) and knowledge base (instances). Also shown are the linkages between text and knowledge base. How such a repository can be used to, for example, enhance the search experience, is illustrated in a subsequent section.

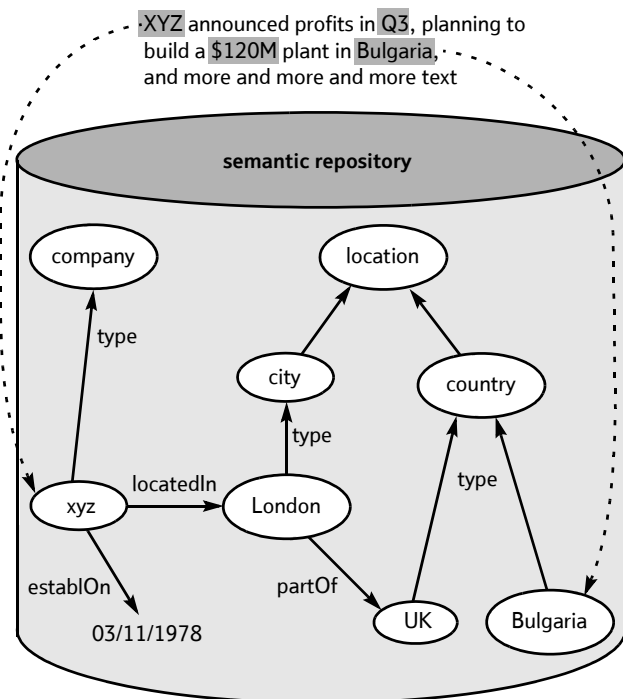


Fig 3 Linking text to knowledge base.
(Courtesy: SIRMA AI EAD [12])

Just as the ontology may consist of an initially manually defined part, plus a learned part, so the knowledge base may consist of instances and property instantiations which are given and others which are subsequently learned. The system should maintain the distinction between the two, since users will wish to distinguish between that which is given, and assumed definitely true, and that which is inferred, and may be erroneous.

4. Ontology evolution

In many applications, once initially developed, ontologies cannot be static but must evolve. Ontology evolution is the timely adaptation of the ontology to changes and the consistent management of these changes. It is not a trivial process, due to the variety of

sources and the consequences of the changes. Therefore, it cannot be performed manually by the knowledge worker. This process is supported by the evolution of the management infrastructure. The first important aspect is the discovery of changes. While in some cases changes to the ontology may be requested explicitly, the actual challenge is to obtain and to examine the non-explicit but available knowledge about the needs of the end users. One way to do this is by analysing various data sources related to the content that is described using the ontology.

A complementary approach is to analyse the end users' behaviour which provides information about their likes, dislikes, preferences or the way they behave. Based on the analysis of this information, changes can be suggested to generate an ontology better suited to the needs of end users. Continuous ontology improvements can be achieved by semi-automatic discovery of such changes, i.e. usage-driven and data-driven ontology evolution.

These two kinds of change discovery are discussed below. Before doing so, we should observe that an important aspect in the evolution process is to guarantee the consistency of the ontology when changes occur, considering the semantics of the ontology change [13].

4.1 Usage-driven ontology changes

In this section we will describe how we can analyse an ontology's usage in order to recommend changes. The usage analysis that leads to the recommendation of changes is a very complex activity. Firstly, it is difficult to find meaningful usage patterns. In a search application, is it useful to discover that many more users are interested in the topic 'industrial project' than in the topic 'research'? Secondly, when a meaningful usage pattern is found, the open issue is how to translate it into a change that leads to the improvement of an application. For example, how do we interpret the information that a lot of users are interested in 'industrial' and 'basic research projects', but none of them are interested in the third type of the projects — 'applied research projects'.

Consider the first example above. If there is no relationship between the concepts 'industrial project' and 'research', then the fact that many more users are interested in the former than the latter is of no use for discovering changes. However, if in the second example, we know that the concepts 'industrial', 'basic research' and 'applied research project' are three subconcepts of the concept 'project', then we could delete the 'unused' concept 'applied research project' or merge it with one of the two other concepts (i.e.

'industrial research' or 'basic research'). Manual effort to do this can be time consuming and error-prone, and the process requires highly skilled personnel, which makes it costly.

The focal point of the approach is the continual adaptation of the ontology to the users' needs. As illustrated above, by analysing the usage data with respect to the ontology, more meaningful changes can be discovered. Moreover, since the content and layout (structure) of an ontology-based application utilise the underlying ontology, by changing the ontology according to the users' needs, the application itself is tailored to the users' needs.

4.2 Data-driven ontology changes

We need to ensure that all ontologies, as well as dependent annotations and meta-data, stay up to date with the document base. One possibility would be a complete re-engineering of the ontology each time the document base changes. But of course, building an ontology for a huge amount of data is a difficult and time-consuming task even if it is supported by tools for automatic or semi-automatic ontology extraction. A much more efficient way would be to adapt the ontology incrementally, i.e. to identify for each change all concepts, instances and properties in the ontology which are affected by this change and to modify the ontology accordingly. Therefore, data-driven change discovery aims at providing methods for automatic or semi-automatic adaptation of an ontology, in accordance with the modifications being applied to the underlying data set.

A number of general prerequisites must be fulfilled by any application which is designed to support data-driven change discovery.

The most important requirement is, of course the need to keep track of all changes to the data. Each change must have associated with it various kinds of information, such as its type, the source from which it has been created, and its target object (e.g. a text document). In order to make the whole system as transparent as possible, not only changes to the data set, but also changes to the ontology, should be logged. Moreover, if ontological changes are caused by changes to the underlying data, the former should be associated with information about the corresponding data modifications.

Optionally, in order to take different user preferences into account, various change strategies could be defined, which permit specifying the degree of influence changes to the data have on the ontology. For example, a user might want the ontology to be updated

because of newly added or modified data. On the other hand, he might want the ontology to remain unchanged if some part of the data set is deleted.

Different kinds of knowledge have to be generated or represented within a change discovery system.

- Generic knowledge

Generic knowledge about the relationship between data and ontology is required, so that the ontology can be changed to allow for newly added or modified data. The implementation of data-driven change discovery methods should be embedded in the context of an ontology extraction system. Such systems, e.g. TextToOnto [14], represent general knowledge about the relationship between an ontology and the underlying data set by means of an ontology learning system.

- Concrete knowledge

Concrete knowledge about the relationship between the data and ontology concepts, instances and properties is needed, because deleting or modifying information in the data set might have an impact on existing entities in the ontology. This impact has to be determined by the application to generate appropriate ontology changes. The concrete knowledge to be stored depends on the way the ontology learning algorithms are implemented.

5. A methodology for building ontologies

Ontology creation will never be an entirely automatic process, and human intervention necessitates a methodology to manage that intervention. This is particularly necessary for decentralised knowledge sharing, such as occurs when organisations come together rapidly to form virtual organisations.

The template of such a decentralised process, which we call DILIGENT, is described here. We cannot yet claim that it is a fully fledged methodology. Here we elaborate on the high-level process, the dominating roles and the functions of DILIGENT.

- Key roles

In DILIGENT there are several experts, with different and complementary skills, involved in collaboratively building the same ontology. In a virtual organisation they are often geographically dispersed. Those who build the ontology may or may not use it. Vice versa, most ontology users will typically not build or modify the given ontology.

- Overall process

An initial ontology is made available and users are free to use it and modify it locally for their own purposes. There is a central board that maintains and assures the quality of the shared core ontology. This central board is also responsible for deciding to do updates to the core ontology. However, updates are mostly based on requests by users working in a decentralised fashion. Therefore the board only loosely controls the process. Due to the changes introduced by users over time, and the on-going integration of changes by the board, the ontology evolves. Let us now survey the DILIGENT process at the next, finer, level of granularity. DILIGENT comprises five main steps — build (1), local adaptation (2), analysis (3), revision (4), local update (5) (see Fig 4).

- Build

The process starts by having domain experts, users, knowledge engineers and ontology engineers build an initial ontology. In contrast to existing ontology engineering methodologies, we do not require completeness of the initial shared ontology with respect to the domain. The team involved in building the initial ontology should be relatively small, in order to more easily find a small and consensual first version of the shared ontology.

- Local adaptation

Once the core ontology is available, users work with it and, in particular, adapt it to their local needs. Typically, they will have their own business requirements and correspondingly evolve their local ontologies (including the common core). In their local environment, they are also free to change the reused core ontology. However, they are not allowed to directly change the core ontology from which other users copy to their local repository. By logging local adaptations (either permanently or at control points), the control board collects change requests to the shared ontology.

- Analysis

The board analyses the local ontologies and the requests and tries to identify similarities in users' ontologies. Since not all the changes introduced or requested by the users will be introduced to the shared core ontology⁴, a crucial activity of the board is deciding which changes are going to be introduced in the next version of the shared ontology. The input from users provides the necessary arguments to underline change requests.

⁴ Please note that it is not the intention of the methodology to merge all user ontologies.

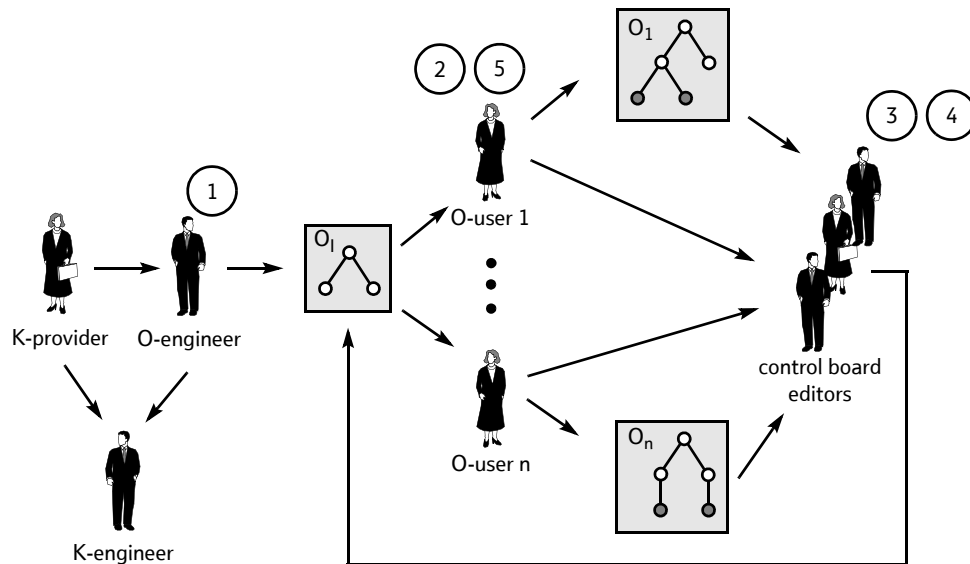


Fig 4 Roles and functions in distributed ontology engineering.

A balanced decision that takes into account the different needs of the users, and meets users' evolving requirements⁵, has to be found.

— Revision

The board should regularly revise the shared ontology, so that local ontologies do not diverge too far from the shared ontology. Therefore, the board should have a well-balanced and representative involvement from the different kinds of participants involved in the process — knowledge providers, domain experts, ontology engineers and users. In this case, users are involved in ontology development, at least through their requests and re-occurring improvements and by evaluating the ontology, mostly from a usability point of view. Knowledge providers in the board are responsible for evaluating the ontology, mostly from a technical and domain point of view. Ontology engineers are one of the major players in the analysis of the arguments for ontology changes and in balancing them from a technical point of view.

Another possible task for the controlling board, that may not always be a requirement, is to assure some compatibility with previous versions. Revision can be regarded as a kind of ontology development guided by a carefully balanced subset of evolving user-driven requirements. Ontology engineers are responsible for updating the ontology, based on the decisions of the board. Revision of the shared ontology entails its evolution.

⁵ This is actually one of the trends in modern software engineering methodologies (see Rational Unified Process [15]).

— Local update

Once a new version of the shared ontology is released, users can update their own local ontologies to better use the knowledge represented in the new version. Even if the differences are small, users may rather use the concepts and properties in the shared ontology rather than analogous locally defined concepts and properties.

It is now widely agreed that ontologies are a core enabler for sophisticated knowledge management systems. The development of ontologies in centralised settings is well studied and established methodologies exist. However, current experiences from projects suggest that ontology engineering should be subject to continuous improvement rather than a one-time action, and that ontologies promise the most benefits in decentralised rather than centralised systems. Hence, a methodology for distributed, loosely controlled and dynamic ontology engineering is needed. Previous work has described a methodology to support the creation of a static ontology in a collaborative ontology engineering setting.

With DILIGENT, we define a process which takes into account that requirements on a knowledge management system change over time. Furthermore, we allow a quick introduction phase with later refinement. Obviously, such a process needs tool support from ontology engineering environments. There exist already some tools which allow for remote and collaborative ontology engineering. However, none exists which could support the complete cycle. We have an implementation which is a first step towards such a tool.

DILIGENT will eventually result in a methodology with tool support to enable ontology engineers to build ontologies in a decentralised environment yet systematically.

6. Accessing knowledge semantically

6.1 Semantic search

The essence of semantic search is to search for a particular semantic entity, not merely a text string as in a conventional search engine. This has a number of advantages.

Firstly, it enables us to disambiguate a text string which has more than one meaning. The person searching for 'Georgia', for example, may be interested in the independent country in South West Asia, or in the state in the USA — or indeed in any one of a number of other entities bearing the same name. With semantic search, users can specify precisely in which Georgia they have an interest. One can go further than this; when searching for a company, for example, one can specify not only that it is a company being sought, but also in which country or industry sector it operates.

Secondly, it enables the identification of an entity even when it has a different representation than the user specifies. Section 3.2 has already described how text entities such as 'George W Bush', 'Mr. Bush', and 'the President' can be semantically equated. Semantic search can make use of this when locating documents. Moreover, a good search engine today will highlight the occurrences of the search string in the documents returned. With semantic search, a search for 'George W Bush'⁶, will also highlight 'Mr Bush', 'the President', or even 'he' or 'him' where appropriate.

Users of search engines do not seek documents, but knowledge. The embedding of semantics in documents means that a search engine can extract the appropriate knowledge from a set of returned documents, and merge that knowledge in a meaningful way, avoiding repetition.

The use of information extraction techniques and the existence of a knowledge base, as discussed above, means that documents returned to the user can be marked up in a highly meaningful way. Colour coding can be used to identify the different classes of entities, e.g. person, company, country. Hyperlinks can be inserted to link the entities to their representation in the knowledge base, enabling the user to be provided with information about the entities; for example, for a

company the user could be informed of its key financial statistics, senior officers, etc.

6.2 Semantic alerting and sharing

Alerting people to the existence of new knowledge, and sharing knowledge between people both depend on the accurate description of users' interests. If profiles are defined too narrowly, users will miss valuable information. If defined too broadly, they will be overwhelmed. Current systems depend upon profiling based on the particular textual language in users' interactions with the system, e.g. the search strings used. With semantic technology, we can go beyond this to define users' interests more precisely with reference to an ontology. The user with an interest in US politics in the 1930s can be made aware of an article on the New Deal, or the influences on President Roosevelt, because the system knows that both New Deal and President Roosevelt are semantically connected to US politics of that time.

The SEKT project (see section 8) has developed a search agent based on Semantic Web technology for alerting users about relevant information. The semantic search agent described allows users to specify semantic queries based upon the PROTON ontology [16]. This returns pages where a named entity (e.g. 'Blair') of a specific type (e.g. Person) has been found. The agent performs a periodic search using the specified semantic query against a set of documents associated with the ontology. The user is then informed (e.g. by e-mail or SMS) when new results are available. The user can then view the results in more detail via a Web interface (see Fig 5).

The ability to allow users to form queries based upon named entities of specific type will improve the precision of the search agent when compared to existing agents that use purely syntactic queries. In addition, the queries specified by the users and their results are being used to enhance the indexing and extraction process of the search engine. The links contained in the pages that form the results of queries form a fruitful source of information related to the query — the pages referenced by the links can be used to seed future indexing and extraction. Users are able to create a number of semantic queries which are executed on a regular basis. Thus, each agent searches for documents that contain entities that match the user's long-term interests. The user is, currently, able to express searches to find documents that contain information about the following:

- a named person holding a particular position, within a certain organisation,

⁶ In practice, we would want to restrict the search, e.g. by specifying a time period or a topic; the former is possible with conventional search techniques, the latter with semantic search.

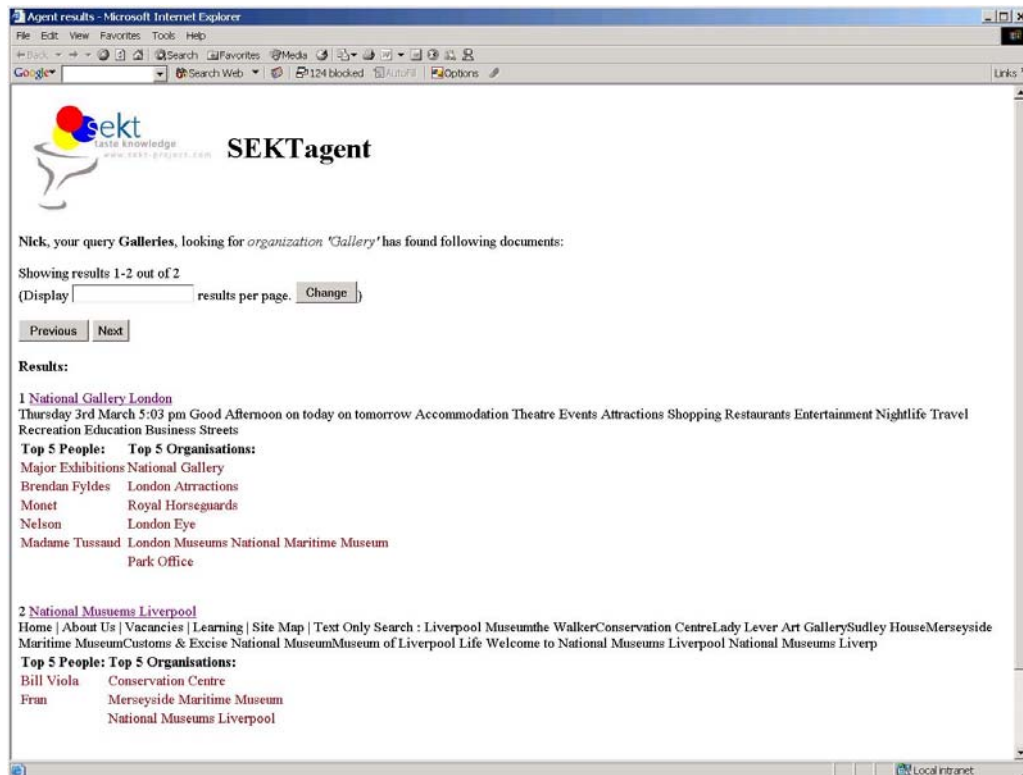


Fig 5 Semantic search agent results.

- a named organisation located at a particular location,
- a particular person,
- a named location,
- a named company, active in a particular industry sector.

Figure 5 shows the results of an agent that was searching for entities of class *Organisation* called 'Gallery'. As can be seen, the result of the query contains the following:

- document title,
- link to the document,
- first 300 characters of the document (taster),
- key people mentioned in the document,
- key organisations mentioned in the document.

The key people and organisations from the document are identified through the ontological annotations linking text in the document to entities in the ontology. They are selected according to number of occurrences, remembering that pronoun resolution and other techniques, described in section 3, will be used to identify all mentions of a given individual ('the Prime Minister', 'Mr Blair', 'he', 'Tony Blair', ...).

Figure 6 then shows how further information can be found from the ontology — by putting the cursor over a given entity from the ontology, the user can obtain a pop-up showing further information about that entity. In the example, the user has requested further information on Merseyside Maritime Museum.



Fig 6 Displaying ontological information.

6.3 Presenting and visualising knowledge

How knowledge is presented, whether as text or in visual form, is crucial to how well and how quickly it can be absorbed. When information about the underlying semantics is available, this can influence presentation. Knowledge represented semantically can be better translated into natural language. Indeed, where knowledge is represented semantically, it can be easily mapped into any number of natural languages. This is not natural language translation, but rather the mapping from structured knowledge (which has been

created by our semantic technologies from unstructured knowledge) into the desired natural language.

Where knowledge is represented semantically, the relationships between elements can be more meaningfully displayed in a visual manner. Here elements may be the entities such as persons, companies, countries described above, or they may be documents in whole or part. The use of semantics to better visualise knowledge is an interesting research area, combining both semantic technologies and the psychology of human-computer interaction.

7. Integrating into applications

The previous section described how the kind of semantic knowledge technology we have been discussing in this paper enables end-user functionality. Examples are:

- searching, browsing and sharing knowledge,
- being alerted to new knowledge,
- knowledge visualisation,
- generating knowledge as natural language.

To be really effective in everyday use, this end-user functionality needs to be integrated into, for example, desktop applications. We can also see these semantic knowledge technologies as part of a value chain, integrated into more specific applications and solutions. A few examples of the possibilities are given in this section. This is by no means an exhaustive list.

7.1 Document and content management

Obvious candidates are applications for document and content management. The former is concerned with the management of large numbers of documents, the latter with the management of items of content, e.g. for the creation of documents. In any case, the capabilities described in previous sections can be valuably integrated into document and content management systems. The ability to semi-automatically create ontologies to describe the material, and to associate items with the ontology, is clearly valuable where large quantities of content need to be managed consistently. The use of semantic search, aided by semantic user profiling, offers improved retrieval capabilities.

We could, for example, envisage our semantic knowledge management applications integrated with the conventional capabilities of a document management system, e.g. configuration management. We could go further and envisage semantic technology intimately supporting the basic document management capabilities. A version control capability which previously listed textual changes between versions of a document could be enhanced to describe semantic

changes, e.g. the inclusion or deletion of an instance, or reference to a relationship.

It is important to understand that we are envisaging more than bringing the basic semantic knowledge capabilities together into the same system as incorporates document management capabilities. This is certainly part of it, but as can be seen from the example above, we can also enhance and add value to the document management capabilities themselves. This is a general point which applies when we integrate semantic knowledge technologies into other applications.

7.2 Information portals

Information portals provide unified access to a range of heterogeneous data sources, in general tailored to the needs of particular users. Much of the functionality discussed under document management could be relevant here. Searching and sharing of knowledge are clearly key to what a user wants from a portal. The use of semantic profiling to enable the sharing of knowledge, without users being overwhelmed by irrelevant information, could also enhance the functionality of a portal.

Knowledge discovery techniques can be used for analysis of usage patterns and tendencies. This can be used not only for user profile construction, but also to detect difficulties with Web site use, to detect portal misuse, and to enhance load-balancing and optimisation. The SEKT project is exploiting usage patterns to develop user profiles, and has developed a SEKTbar to display an ontology of a user's interests, as shown in Fig 7.

A principal role of a portal is in overcoming heterogeneity; this means that the ability to merge ontologies and map between ontologies come to the fore. By mapping between the user's ontology, and that of the individual knowledge sources, we can provide one integrated view tuned precisely to the user's requirements.

Some portals are designed to provide an interface to specific processes, in general the processes by which an organisation is managed. The use of semantic technology to describe, locate and create processes is explained by Davies et al [17], while Duke et al [18] describe a specific application of this approach. With this approach a portal can be used to rapidly create a new process and to manage existing processes.

7.3 Integrated business communications

Greater connectivity between personal software tools such as diary and e-mail software would have big



Fig 7 Screenshot, with SEKTbar on the left. At the top, the SEKTbar shows the ontology of the user's interests. Below that, is shown the list of Web pages which correspond to the semantic-web-rdf topic.

benefits for personal productivity. Consider the linking of diary and e-mail software. This could aid prioritisation of e-mails. Knowing that the diary contained an imminent meeting with a particular person would enable the system to prioritise any e-mails from or about that particular person or his or her organisation. Hyperlinks could be created from the diary entry to all the relevant e-mails. Diary and e-mail entries could also be linked to information in the corporate intranet, or Web. This could be done directly. Alternatively a semantic repository comprising ontology and knowledgebase could be constructed from the relevant

information, and linkages established to the knowledge base. These possibilities are illustrated in Fig 8.

By using machine reasoning, connections can be established which would not immediately be apparent. To take a trivial example, a diary entry might refer to company A, which is in fact a subsidiary of company B. Knowing this fact, the system could present the user with breaking news about company B prior to the meeting with company A.

This functionality can be extended to a group of collaborating individuals, where linkages between different individuals' personal productivity software could be valuable. An example of this would be creating a linkage from one person's diary entry referring to a particular company or person, and another person's e-mails about or from that particular company or person. This is shown in Fig 9. The approach could be extended to a wide range of personal software tools.

7.4 Business intelligence and customer relationship management

Until recently, the term 'business intelligence' has referred to the analysis of an organisation's structured data, e.g. held in relational databases. Typical applications are customer profiling for cross-selling and understanding product trends for inventory management.

It is now realised that a great deal of an organisation's knowledge resides as unstructured textual data, and business intelligence is being enhanced to take account of such data and to merge the knowledge so gained with what can be learned from structured data. Besides an organisation's own data, other sources of unstructured text can be mined, e.g. the Web, analysts' reports, and competitors' literature. Examples are:

- analysis of customer e-mails and customer input to the company Web site — this could be used to detect and categorise customer concerns and issues,

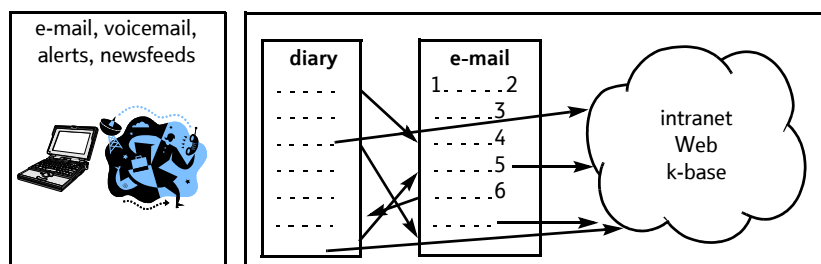


Fig 8 Applying Semantic Web technology to personal productivity tools.

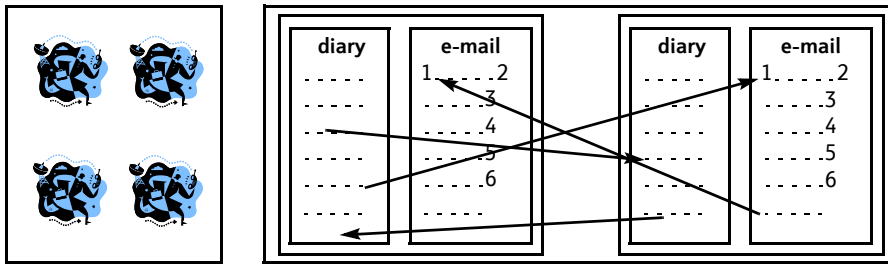


Fig 9 Linking team-members' personal productivity tools.

- analysis of internal company information, e.g. e-mails, to detect generic company issues as they arise — analysis of information created by salespeople could detect sales issues and changes in market conditions, or help understand employee concerns, the last of which raises an issue that is often present when using this technology, that of confidentiality (monitoring employee e-mails to detect employee concerns could only be done to produce generic information, not information specific to individuals),
- analysis of competitor information, ranging from sales literature to patent information, to detect trends in products and marketing,
- 'sentiment analysis' on the Web, i.e. crawling the Web to detect reaction to one's own products, or to competitors.

Closely related to business intelligence is customer relationship management. The difference being that, while the former is concerned with trend analysis, including relating to customer habits and preferences, the latter is aimed at collecting information about specific customers. This information can then be used to target customers appropriately. The basic principle is the same, though. By analysing corporate knowledge, in the form of unstructured text, we can augment what is already known from structured knowledge.

In addition, we can analyse how customers navigate the company Web site so as to gain commercially valuable information about their habits and preferences. This is quite apart from the use of semantic knowledge applications to enhance the customer's experience when visiting the Web site, and to aid call-centre staff when dealing with customers.

7.5 eLearning

Our final example is different from the preceding in that, while it builds on them, it is targeted at a specific sector, i.e. education and training. There is a strong link here with content management. The requirement is to manage knowledge elements and combine them to form an eLearning experience.

Learning object management systems have been developed to provide infrastructures for universities. On a smaller scale they are used by publishers to manage a collection of learning materials, in order to permit more targeted and flexible reuse. To support this, the use of ontologies is being explored, e.g. by the Learning Lab centred at Hannover [19].

8. The SEKT solution

BT and its partners are actively developing next generation knowledge management solutions based on semantic knowledge technology. BT is leading a collaborative project, SEKT [20], which is combining the technologies described in this paper to create next generation knowledge management solutions.

In particular, SEKT is using:

- the GATE architecture for natural language processing, developed at the University of Sheffield's natural language processing research group [11],
- the Text-Garden text mining software [21] from the Jozef Stefan Institute in Ljubljana,
- ontology management software from the AIFB Institute at the University of Karlsruhe [22].

Apart from leading the project, and developing the intelligent content management case study described below, BT is leading the development of the knowledge access tools, e.g. search and browse, available to the end users.

SEKT is also drawing on the knowledge management solutions expertise of Empolis [23], and research and ontology engineering expertise from a number of other prominent European research teams.

8.1 The SEKT integration platform

The various components of SEKT need to be tightly integrated. One could, in principle, define each component as a Web Service and use the SOAP protocol

to achieve integration. However, for commercially scalable systems, a different approach is required.

This approach, which uses an architecture developed by Empolis, is based on combining computational objects called pipelets to form pipelines. A pipelet represents a particular service performed by a SEKT component. Pipelines consist of pipelets with a flow of control. As well as consisting of pipelets in sequence, a pipeline can incorporate conditional branching, and also one pipeline can call another, in a manner analogous to a subroutine call.

Pipelets communicate via a central data structure called a 'sync board'. By avoiding the need for pipelets to communicate directly, this not only keeps pipelets independent but also minimises communications overhead. This is because data structures are only passed to pipelets if they request them. The data representation on the sync board will comply with a W3C standard, specifically OWL-Lite. This is the simplest of the three variants of the OWL language [24].

8.2 *An open architecture*

Integration of the SEKT technologies is necessary, but not sufficient, for next generation knowledge management. It is also necessary to be able to integrate SEKT software into knowledge management solutions. For that, SEKT is developing APIs which will be published to encourage third party developers to take the SEKT infrastructure and build upon it for their own applications. Some aspects of these APIs are still a matter of research, in particular APIs for ontology management. However, the development of these APIs will be consistent with emerging standards.

8.3 *The SEKT case studies*

Within SEKT there are three case studies. Their role is both to test out the technical feasibility of our approach, including its scalability, and also to understand what functionality users need, and how they want to interact with that functionality.

The largest case study is within the Spanish legal system. Here, newly qualified judges, 'in their first appointment', are confronted with difficult and significant decisions in situations where they need to make these decisions on their own. The solution consists of two parts, which must be integrated together at the semantic level. Firstly, a database of frequently asked questions is being developed and semantic technology is being used to match a judge's question to the most appropriate answer. However, this is not in itself sufficient, since the judge must also be able to defend their decision by reference to the appropriate articles under Spanish law. There are

various legal databases which are relevant here, each with their own ontology. Ontology mapping will be needed between the ontology of the question and answers database and that of the legislative databases. Ontology mapping techniques being developed within SEKT will enable this.

Another case study comes from knowledge management. SEKT is developing an application to serve the needs of several thousands of IT consultants, distributed across the globe. Each consultant has their own way of conceptualising their view of their domain. The application will enable knowledge to be shared, and yet viewed by each consultant using an ontology which is natural to them.

The third case study uses semantic technology for intelligent content management, specifically within the BT digital library. Here SEKT technology will provide users with more precise querying and browsing capability. In addition, a digital library is a platform for sharing knowledge, and SEKT's semantic technology will be used to target knowledge sharing more precisely. Another goal of this case study is to provide a common view to knowledge from a wide range of sources. This might be based on an ontology shared by all users, or it might be through an ontology created for each user, and tuned to that user's requirements. In any case, focused crawling is being used to gather knowledge from the corporate intranet, or even potentially the Web, to augment knowledge already held within the digital library. Moreover, we know that much valuable knowledge within a corporation is stored on the desktop. As far as is possible given the constraints of privacy and confidentiality, this will also be being brought into the common framework of the digital library.

8.4 *The human dimension*

The kind of knowledge management systems we have been discussing in this paper are only possible by using the most advanced technology. However, they will only be used effectively if we understand how people interact with such systems. This dimension is not being neglected in SEKT. One of the SEKT partners, Kea-pro [25], provides specialised consultancy in IT usability. Each of the case studies will undertake a programme of usability validation and the results from this will be fed back into the technical development. In this way, the final system will not only be demonstrated to be technically feasible, but also fit for purpose as a tool for knowledge workers.

9. **Conclusions**

This paper has described how knowledge management will be improved through an understanding of the underlying semantics of information. To realise this, a

number of different technologies need to be brought together. These technologies pre-date the ideas of the Semantic Web. However, their fusion provides the infrastructure which makes the Semantic Web and semantic knowledge management possible. There is still a great deal to do. A major challenge is the development of semi-automatic techniques for ontology mapping; mapping between ontologies is at the heart of overcoming heterogeneity in semantic knowledge systems. There is also work to be done to understand how these underlying technologies can optimally work together, and how they can be integrated into applications. There is also work to be done at the human level to understand how people can best use the technology, and to understand how it can really help knowledge workers. However, there are already applications of the technology being used to solve real business problems, and the next few years will see many more.

Acknowledgements

The work described in this paper was undertaken as part of the SEKT project, funded by the European Commission under the 6th Framework Programme (IST-2003-506826). The authors would like to acknowledge the input of their colleagues in the project. Section 2, in particular, is based on reviews on the state of the art in knowledge discovery and information extraction, undertaken respectively by the Jozef Stefan Institute [21], Slovenia, and the University of Sheffield [11]. The SEKT Integration Platform described in section 8 is being developed by Empolis [23].

References

- 1 Antoniou G and van Harmelen F: 'A Semantic Web Primer', The MIT Press (2004).
- 2 Davies N J: 'Knowledge management', *BT Technol J*, **18**, No 1, pp 62—63 (January 2000).
- 3 Maxwell C: 'The future of work — understanding the role of technology', *BT Technol J*, **18**, No 1, pp 55—56 (January 2000).
- 4 Lenat D B and Guha R V: 'Building large knowledge-based systems: representation and inference in the Cyc project', Addison-Wesley (1990).
- 5 dmoz project — <http://dmoz.org/>
- 6 Agirre E, Ansa O, Hovy E and Martinez D: 'Enriching very large ontologies using the WWW', in Proceedings of the First Workshop on Ontology Learning OL-2000, the 14th European Conference on Artificial Intelligence, ECAI-2000 (2000).
- 7 Steinbach M, Karypis G and Kumar V: 'A comparison of document clustering techniques', in Grobelnik M, Mladenic D and Milic-Frayling N (Eds): 'Workshop on Text Mining', Proc KDD, pp 109—110, Boston, MA, USA (2000).
- 8 Mitchell T M: 'Machine Learning', The McGraw-Hill Companies Inc (1997).
- 9 Hotho A, Staab S and Stumme G: 'Explaining text clustering results using semantic structures', in Proceedings of ECML/PKDD 2003, LNAI 2838, pp 217—228, Springer Verlag (2003).
- 10 Mladenic D and Grobelnik M: 'Mapping documents onto web page ontology', in Berendt B, Hotho A, Mladenic D, Someren M W, van Spiliopoulou M and Stumme G (Eds): 'Web mining: from web to semantic web', Lecture notes in artificial intelligence, LNCS 3209, pp 77—96, Springer (2004).
- 11 University of Sheffield — <http://nlp.shef.ac.uk/>
- 12 Sirma — <http://www.sirma.com/>
- 13 Stojanovic L, Maedche A, Motik B and Stojanovic N: 'User-driven ontology evolution management', in European Conference on Knowledge Engineering and Management (EKAW 2002), pp 285—300, Springer-Verlag (October 2002).
- 14 TestToOnto — <http://www.sourceforge.net/projects/texttoonto/>
- 15 Rational Unified Process — <http://www-306.ibm.com/software/awdtools/rup/index.html>
- 16 PROTON — <http://proton.semanticweb.org/>
- 17 Davies N J, Fensel D and Richardson M: 'The future of Web Services', *BT Technol J*, **22**, No 1, pp 118—130 (January 2004).
- 18 Duke A, Davies N J and Richardson M: 'Enabling a scalable service oriented architecture with semantic Web Services', *BT Technol J*, **23**, No 3, pp 191—201 (July 2005).
- 19 Learning Lab — <http://www.learninglab.de/>
- 20 SEKT — <http://www.sekt-project.com/>
- 21 Text Garden — <http://textgarden.ijs.si/>
- 22 AIFB — <http://www.aifb.uni-karlsruhe.de/>
- 23 Empolis — <http://www.empolis.com/>
- 24 OWL — <http://www.w3.org/2004/OWL/>
- 25 Kea-pro — <http://www.keapro.net/>



Dr John Davies leads BT's Next Generation Web Research Group at Adastral Park. Current interests centre around the application of Semantic Web technology to knowledge management, information retrieval and semantic Web Services. He is industrial chair of the Semantic Web Services Initiative, co-organiser of the European Semantic Web Conference series and Project Director of the SEKT EU integrated project. He has written and edited many papers and books in the areas of web-based information management, knowledge management, virtual communities and the Semantic Web and has served on the program committee of many conferences in related areas. He is a member of the British Computer Society and a Chartered Engineer. Earlier research at BT led to the development of a set of knowledge management tools which are the subject of a number of patents. He founded and served as CTO with Exago Ltd which was set up to exploit this technology. Exago was subsequently acquired by Corpora Ltd, with whom he retains the role of Group Technical Adviser.



Rudi Studer is Full Professor in Applied Informatics at the University of Karlsruhe, Institute AIFB. His research interests include knowledge management, Semantic Web technologies and applications, ontology management, data and text mining, Web Services, and peer-to-peer systems.

He is one of the presidents of the FZI Research Center for Information Technologies at the University of Karlsruhe as well as co-founder of the spin-off company Ontoprise GmbH that develops semantic applications. He is engaged in various national and international cooperation projects, among others the EU funded Integrated Projects SEKT or the Graduate School IME. He is president of the Semantic Web Science Association and one of the Editors-in-chief of the journal *Web Semantics: Science, Services, and Agents on the World Wide Web*.

Next generation knowledge management



York Sure is an assistant professor at the Institute AIFB (University of Karlsruhe). He graduated in Industrial Engineering and received in 2004 his PhD in Computer Science from University of Karlsruhe. He published over 40 research papers in the areas of ontology management, semantic web and knowledge management. He previously worked for the EU IST On-To-Knowledge project and is currently leading the work of the AIFB in the EU integrated project SEKT and the EU thematic network Knowledge Web. He is a member of IEEE, ACM and the German associations for

Informatics (GI) and Knowledge Management (GfWM). He serves as a reviewer for international journals, conferences, and workshops. He is co-inventor of the German national conference series on knowledge management 'Wissensmanagement — Erfahrungen und Visionen' (held in 2001, 2003, 2005). For the European Semantic Web Conference (ESWC 2006) he has been nominated as program chair.



Paul Warren works in BT's Next Generation Web research group, where he is SEKT project manager and also responsible for the project's exploitation strategy.

In his previous role within BT he undertook technology foresight, advising BT of the likely impact of future technologies. It was in this role that he became aware of the commercial and industrial significance of the Semantic Web.

He has published widely on technology management, technology foresight, and recently the application of the Semantic Web.

His original degree was in theoretical physics, and he subsequently gained an MSc and MPhil in electronics.