

Dynamic Data Fusion Using An Ontology-Based Software Agent System

Mark T. Elmore Thomas E. Potok and Frederick T. Sheldon¹

Applied Software Engineering Research Group
Computational Sciences and Engineering Division
Oak Ridge National Laboratory²

Abstract

Developing a knowledge-sharing capability across distributed heterogeneous data sources remains a significant challenge. Ontology-based approaches to this problem show promise by resolving heterogeneity, if the participating data owners agree to use a common ontology (i.e., a set of common attributes). Such common ontologies offer the capability to work with distributed data as if it were located in a central repository. This knowledge sharing may be achieved by determining the *intersection* of similar concepts from across various heterogeneous systems. However, if information is sought from a subset of the participating data sources, there may be concepts common to the subset that are not included in the full common ontology, and therefore are unavailable for knowledge sharing. One way to solve this problem is to construct a series of ontologies, one for each possible combination of data sources. In this way, no concepts are lost, but the number of possible subsets is prohibitively large.

We offer a novel software agent approach as an alternative that provides a flexible and dynamic fusion of data across any combination of the participating heterogeneous data sources to maximize knowledge sharing. The software agents generate the largest intersection of shared data across any selected subset of data sources. This ontology-based agent approach maximizes knowledge sharing by dynamically generating common ontologies over the data sources of interest.

The approach was validated using data provided by five (disparate) national laboratories by defining a local ontology for each laboratory (i.e., data source). In this experiment, the ontologies are used to specify how to format the data using XML to make it suitable for query. Consequently, software agents are empowered to provide the ability to dynamically form local ontologies from the data sources. In this way, the cost of developing these ontologies is reduced while providing the broadest possible access to available data sources.

Introduction

The Department of Energy (DOE) manages the United State's national laboratory system, each of which has evolved a variety of business models for managing research proposals over the past six decades. Because of the historical nature of these evolutions, both the business models, and their associated (heterogeneous) data collections, are deeply rooted. A system was needed that could merge data from the heterogeneous systems as if the data was gathered and stored in a centralized repository.

For example (as shown in Figure 1 upper left hand corner), each lab within DOE has a slightly different way to collect research proposal data. Each collection and its corresponding ontology is based on a unique business model. Moving to a single common ontology to resolve data heterogeneity across all laboratories will result in an intersection of

all data concepts. Unfortunately, only common concepts are retained. Moreover, if two or more labs have a concept that is not common, then (potentially) valuable information is not being represented. In this case, a common ontology that results in loss of such data has an unacceptable impact on the associated business models.

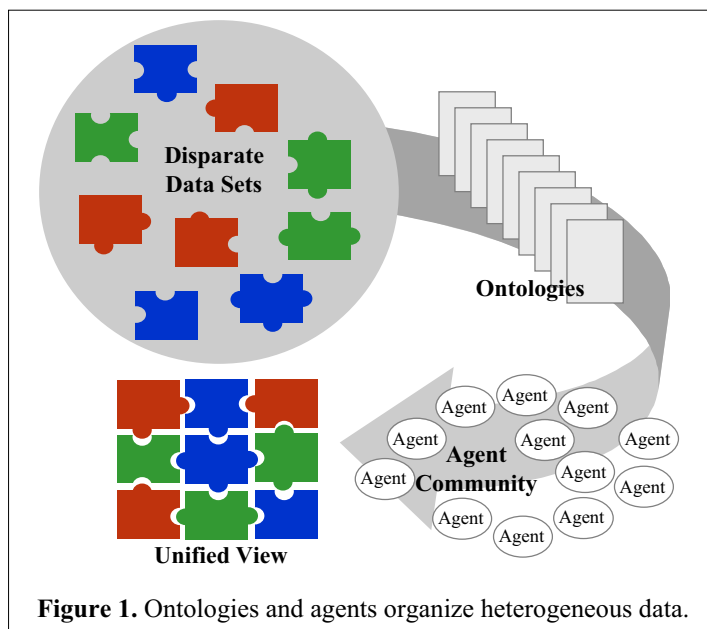


Figure 1. Ontologies and agents organize heterogeneous data.

¹ Contact: Post Office Box 2008, MS 6363, Oak Ridge, TN 37831-6363, Phone: 865-576-1339, Fax: 865-574-6275, SheldonFT@ornl.gov.

² The submitted manuscript has been authored by a contractor of the U.S. Government under contract No. DE-AC05-96OR22464. Accordingly, the U.S. Government retains a non-exclusive, royalty-free license to publish or reproduce the published form of this contribution, or allow others to do so, for U.S. Government Purposes.

Figure 2 gives an example of how concepts can be lost in a common ontology. Consider three data sets labeled *A*, *B*, and *C* each of which contain three concepts, *Lab Name*, *PI Name*, and one other that is not common. Lets examine the merged data from data sets A and B. The first two concepts match across all three but the third concept is lost. This is unfortunate because we are concerned only about data from sets A and B (see Figure 2 right lower corner). A better method is needed that would alleviate this unfortunate problem.

In this case, the new approach needs to add *Duration* to the ontology when only data sets A and B are being considered. If this is possible, then full visibility of all data across all combinations all participating data sets could be accomplished with a *series* of merged ontologies.

The ontologies are merged in a way that provides for each possible combination of source data concepts, as compared to a single common ontology. On the other hand, there are a large number of possible combinations a user may choose potentially generating a huge number of ontology combinations. Current ontological approaches for merging heterogeneous data sets have been successful, but require *all* data owners to participate in building a single common view. For this reason, a large number of possible data set combinations ontologies for every possibility cannot be “pre-built.” Instead, we propose to use software agents to build the desired ontology combination on-the-fly for each user-generated merge and query/search operation (as shown in Figure 1 bottom half). In this way only relevant merges are implemented, avoiding the need to generate all possible combinations while satisfying the possibility for merges that consider all possible combinations (including the intersection of every data set).

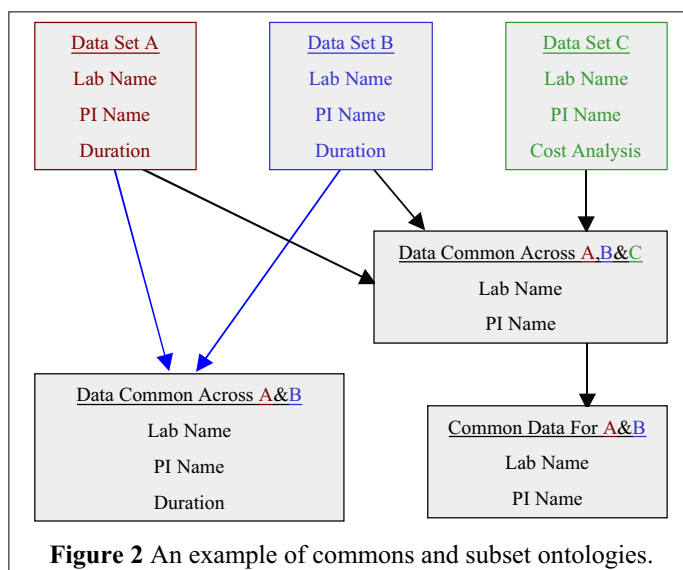


Figure 2 An example of commons and subset ontologies.

Background

The term “ontology” was first used in philosophy (Wolffe, C., *Philosophia Prima Sive Ontologia*, 1729) where its use means “a branch of metaphysics relating to the nature and relations of being.” Computational scientist’s adopted the general meaning of the term to be the “relations of being” for describing how data is related. Ontologies have proven to be a useful tool for data integration across heterogeneous data sets.

In effect, an ontology specifies a vocabulary that is used to discuss a problem domain. For example, an ontology for baseball would include terms such as ball, bat, glove, strike, foul, etc. But it is not so much what terms are *used*, but what those terms *mean* that reaches to the core of how ontologies are used. For example, changing the language of an ontology from English to French changes the terms used, but does not change the concepts specified by the terms [1]. Gruber sums this up well when he describes an ontology as an “explicit specification of a conceptualization”[2].

In several revisions from 1993 to 1995, Gruber [2] represents a foundational set of design criteria to guide the development of ontologies in support of knowledge sharing activities. Gruber applies formal engineering discipline to ontology design using a core set of five design criteria as follows. The ontology should provide *clarity* in defining terms and *coherence* by being logically consistent. It should provide *extendibility* to allow expansion without affecting existing definitions, and should have *minimal encoding bias* so that the notation used to describe a concept does not restrict alternative ways to understand the concept. Finally, an ontology should have *Minimal ontological commitment*, meaning that the description of concepts should be as loose as possible to permit flexible use of the described concepts. These seminal design criteria form the basis for a series of ontological studies.

Recent in this series, Holsapple and Joshi [ref here] adopt Gruber’s view and give a taxonomy of five approaches to Gruber’s ontology design criteria. First lets consider the *inspirational approach*. In this case, an individual uses their own insights and viewpoint to develop an ontology of the domain that is then (hopefully) adopted by other users. The *inductive approach* builds an ontology based on a specific (set of) case(s) within the domain and is refined by evolving toward a more generalized ontology. The *seductive approach* moves in the opposite direction, beginning with an ontology built upon general principles of the domain and evolving toward fulfilling specific cases. The *synthetic approach* is a combination of existing ontologies into a single all-encompassing ontology that describes the combined domain. Finally, in the *collaborative approach*, a team of individuals incorporates aspects of the other approaches to build an ontology using the combined viewpoints and

(possibly) using existing ontologies as an anchor. Holsapple and Joshi find the collaborative approach the most useful for their problem, and present a case study using it.

These various ontology based approaches, from Gruber to Holsapple and Joshi, provide methods to discover the homogeneity that may be found among heterogeneous data sets, and from that, build a common ontology. However, these approaches assume that participants are capable of migrating their data to a new ontology; this was not the case in our problem. The varied data collections within DOE are tightly coupled to their associated business processes. So much so, that it became unclear whether the business model drives the data or the data drives the business model. Migrating the data to a common ontology would necessitate a prohibitively expensive change to long and well-established business models. Thus, these powerful approaches are not suited to our problem without a variation from previously seen ontology-based data fusion.

Our approach

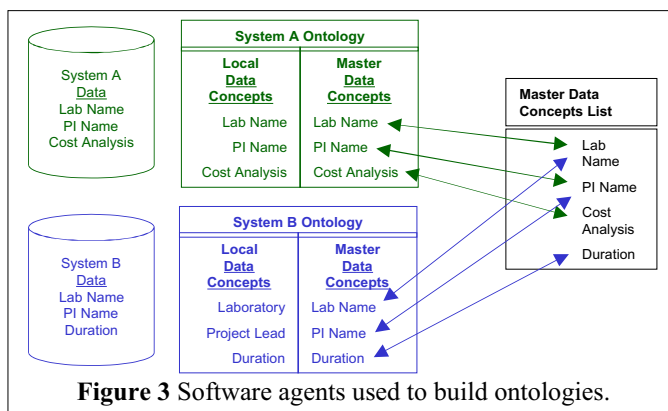
Given the aforementioned ontological approaches to data merging we decided that the power of these tools showed great promise. Current ontological approaches to the merging of heterogeneous data have been successful, but require the owners of the data to participate in the adoption of a single common ontology. In our case, we cannot go to a single, all encompassing, common ontology because of the importance of concepts occurring in most, but not all data sources. The solution to this would be to build a series of ontologies for all possible combinations of the underlying data sets. However, in our case this is an intractable problem and we decided to explore the possibility of using software agents to perform the ontology building tasks automatically.

Foundational to our approach is the use of the Extensible Markup Language (XML) as the mechanism for capturing data [3]. Each laboratory has a different mechanism for capturing data, from databases, to spreadsheets, to ad hoc text files and various combinations. Using a series of specialized software tools, the various data formats are converted into XML. XML provided an efficient mechanism to bring the distributed and heterogeneous data formats into a powerful, flexible, and common format thereby providing a standard ontology-to-software agent interface [4]. Consequently abstracting unnecessary details from the underlying datasets, which proved to be a sufficiently rich environment for the software agents to perform the merging of ontologies.

Next, we designed the coordinator for the system, a mechanism that would allow software agents to understand the various data sets and enable the agents to merge ontologies. While formal ontologies have a great many strengths, one potent drawback is the learning curve associated with using them. In our problem, the data owners who understood the data did not have the time for such a learning curve. To address this conflict, we devised a simple system that specifies data concepts by defining the equivalence relations between a data owner's local data concepts and other participating data concepts. Our technique works in the same way that the concept for "table" can be specified by showing equivalence to an already understood concept of the Spanish word "mesa," and/or the French word "tableau," and/or the Italian word "tavolo," Although intentionally informal in its execution, the specification of the data concepts by their relationships meets Gruber's definition of an ontology.

For example, as shown in Figure 3, the first data set, *Data Owner A* provides a simple list of data elements. In a parallel column the list of elements is repeated. The first column is labeled *Local Data Concept List*; the second column is labeled *Master Data Concept List*. Together they represent a simple ontology, i.e., a specification of the data concepts represented by a mapping from the *Local Data Concept List* to the *Master Data Concept List*. Next we incorporate data from *Data Owner B* where data owners A and B know each other's data. This is the case across most of the DOE system; most data owners know their data as well as its relationship to a small number of the other data systems. This knowledge context provided the basis for constructing composite ontologies (i.e., master data concept).

For the second data set, *Data Owner B* looked at *Data Owner A's* ontology mapping. *Data Owner B* could then provide us the knowledge needed to build a new ontology mapping, mapping B's data to the *Master Data Concept List*. Each data concept in a local data list was mapped to the same concept in the *Master Data Concept List*. If *Data Owner B* had data that was not in the *Master Data Concept List*, a new entry was added to the *Master Data Concept List*. Conversely, if a data concept in the *Master Data Concept List* was not present in the local ontology then there was no



mapping established from his *Local Data Concept List* to the *Master Data Concept List*. In the end, the system consists of three elements: (1) *Master Data Concept List*, (2) *Local Data Concept List to Master Data Concept List Ontology* for Data Set A, and (3) *Local Data Concept List to Master Data Concept List Ontology* for Data Set B.

This process continues for each new data set. Each new data owner uses the previous work to help determine their ontology as a specification of the mapping from their *Local Data Concepts List* to the *Master Data Concepts List*. Questions pertaining to proper mappings are resolved (through discussion) between the new data owner and the data owner that previously added the data concept to the *Master Data Concepts List*.

Thus, in this process, the *Master Data Concepts List* is a union of the data concepts across all participating data sets, and a given data set's ontology is a mapping specifying the relationships between the intersection of that data set's local data concepts and the master data concepts. Relationships among a selection of the local data sets' ontologies can be determined using the *Master Data Concepts List* as a point of common reference. It is interesting to note that there is no centralized ontology for the entire system. Instead, it is distributed across the ontology mappings of the individual data sets and the *Master Data Concepts List*. For example, *Data Set A's* "PI Name" specifies the same concept as *Data Set B's* "Project Lead" but this cannot be directly determined at one centralized point; rather it is determined via the data set ontologies and the *Master Data Concepts List*. Software agents use this distributed ontology to provide the functionality of a centralized ontology along with the ability to be flexible in meeting the varied needs of the users.

As described above, we applied this approach on data from five national laboratories. These laboratories are very large and present massive data sets across a diverse set of repositories (e.g., databases, spreadsheets, and simple ASCII files, etc.). We manually created the local ontologies, and used these ontologies to create XML representations of the data at the laboratories. From this base, we then applied our approach to building dynamic on-the-fly ontologies using agents.

Agents and ontologies

For each data set, we establish a *Data Agent* (i.e., a software agent) that is assigned the responsibility of knowing how to retrieve data from the underlying data repository as well as present the data to the overall system for merging. The data agent uses the ontology that was built by the data owner for that local data set. That ontology mapping allows the data agent to act in a bilingual manner, that is, to understanding both the local and master concepts. Based on the master data concepts side of the ontology, the agent understands the language of the agent community; based on the local data concepts side of the ontology, he understands how to retrieve his local data; and based on the mapping between the two sides of the ontology, he understands how to translate between the two.

A *Data Integration Agent* is the final piece of the system that builds merged ontologies, on-the-fly, according to a user's desires. The Data Integration Agent directs and coordinates the activities of the data agents in the system's software agent community. The Data Integration Agent is also responsible for accepting requests from a *Graphical User Interface (GUI) Agent*, another important type of bilingual agent, responsible for translating user requests into the language of the software agent community, and then translating the results into a visualization useful to the user. The Data Integration Agent distributes the GUI Agent's requests to the appropriate data agents, merges the results from the data agents, and passes the merged data back to the GUI Agent for visualization as shown in Figure 4.

Because this is a software agent community, the modules of the system, in the form of software agents, are very loosely coupled. This is true (first) because each agent can be built in any manner, irrespective of the other agents but compliant to the agent community's web-based interface. Second, the agents may be geographically separated, operating from any accessible point on the web. Third, there is fault tolerance in the system; the system will continue to function if some data agents are unavailable. A missing agent obviously cannot contribute data to a solution, but will not prevent a solution from being created. Fourth, there may be multiple GUI agents to meet a variety of preferences from the user community. Clearly, the agents and ontology approach provides an eloquent efficient and extensible solution.

Software agents build a common ontology

When a user of the system first brings up the GUI, the GUI Agent asks the Data Integration Agent for a list of available data sets, that is, a list of available data agents. The Data Integration Agent then checks the list of data agents that are registered, and verifies the availability of each. This agent then reports the availability to the GUI Agent, who displays the available data sources to the user. The user then selects the desired data sources and the software agents dynamically create a merged ontology for the selected data sources.

To create this merged ontology, the Data Integration Agent sequentially distributes the Master Data Concepts List to the data agents chosen by the user. The first data agent compares this concept list to his local ontology, and deletes from this list the data concepts that are not found (i.e., the data concepts that are not in the local ontology). The data agents then hands the reduced data concept list back to the Data Integration Agent who then passes the reduced list to the next data agent selected by the user. This process continues for each of the user-selected agents until all have seen the list.

The trimmed data concepts list resulting from this process is the intersection of the data concepts captured within the participating systems. Note how this process is the reverse of the way the Master Data Concepts List was originally generated. In building the Master Data Concepts List, the data owners each *added* concepts from their local data sets that are new concepts to the Master Data Concepts List. Here, the data agents *remove* data concepts that are not part of their local ontology from a copy of the Master Data Concepts List

The final reduced data concepts list, in conjunction with each participating data agent's ontology, constitutes a shared ontology across the participating data sets. This ontology is dynamically generated based on a request from a user, and is evaluated against the latest information from each local data source. Participating agents can each understand and provide information about all the data concepts that are shared across the participating systems, which significantly increase the capability of current ontologies. Moreover, the full system addresses Gruber's five ontology design criteria. Using relationships to specify concepts provides *Clarity* and *Coherence*. *Extendibility* has been shown in the process by adding new data sets. And the succinct specification renders *Minimal Encoding Bias* and *Minimal Ontological Commitment*. This process could also be thought of as the software agent version of Holsapple and Joshi's *collaborative approach* to ontology design, producing results in much the same way a group of collaborating humans would have done, but significantly faster, and with far greater accuracy.

Querying over merged data

The GUI offers the human user an interface to specify a query over distributed data as if it were collected in one location under a single schema. Our version of the GUI was designed to provide the user with rich query capabilities to permit users to specify down to the data fields of interest (as in an SQL *select*) and, at the same time, permitting constraints on data field values (as in an SQL *where*).

Queries are passed to the Data Integration Agent who partitions the query out to the selected set of data agents. Agents are then responsible to fulfill the query over their local data set. The underlying XML data set may be stored in any manner usable by the data agent. Each data agent uses the local ontology mapping to convert the query request from the Data Integration Agent into a format usable against the respective data set. The answer to the query is then translated back into the format used by the software agent community and returned to the Data Integration Agent. For example, from Figure 3 a general search for a PI Name of "smith" would be translated to a search over a Project Lead with the name "smith" within System B. Both the query and the response are XML format. The XML format enables the Data Integration Agent to assemble the individual responses from the data agents. The data integration agent then passes the assembled query response back to the GUI agent for presentation to the user.

Results

After completing the system design, five data owners were contacted and asked for their assistance by participating in a prototype demonstration. This selection was of sufficient size to test our approach; five data sets producing ($2^5 - 1$) or 31 possible combinations of merged data sets over which a user might wish to query. This number is too large for the data owners to build all 31 possible ontologies, yet large enough to test the strength of our software agent enabled ontology building approach.

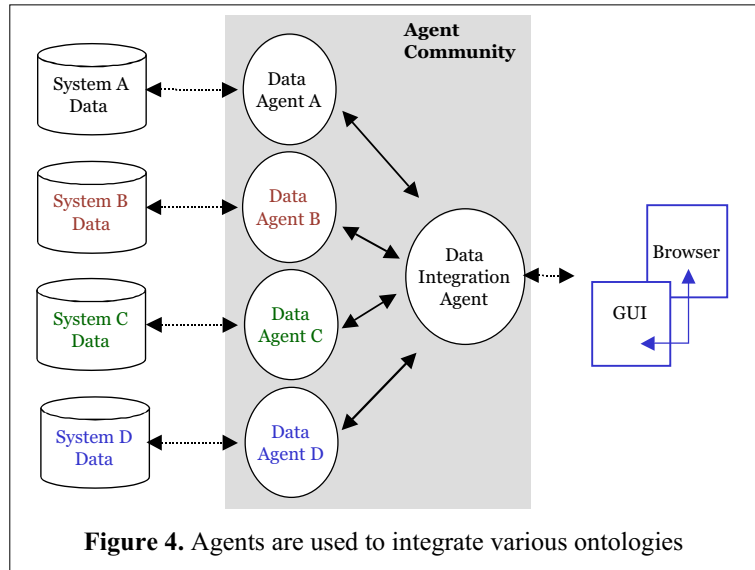


Figure 4. Agents are used to integrate various ontologies

The data owners were enthusiastic about the potential solution yet skeptical about using software agents as described above. In fact, all prior efforts with the same goal were unsuccessful. These data owners are extremely busy; they could only give small amounts of time to help validate our approach. Indeed, the approach proceeded better than expected. It took more time to explain the approach to the data owners than to actually incorporate their data. Some data took longer than others to incorporate (i.e., define in terms of a compliant ontology) because varied native formats, but none of the data took more than a couple of days.

The system performed without a glitch. Acquisition times for query results were negligible, with network latency being the bottleneck. Delays were similar to downloading a web page of typical complexity, well within most typical users' tolerances for delay. The sponsor deemed the prototype a success. Due to the success of the prototype, the DOE decided to implement the system across all of its installations, estimating a cost savings of \$39 million per year. Consequently, the DOE has now begun the process of building a production system based on our prototype.

Conclusions and future work

Ontology based merging of data has proven to be a viable technology in a great many instances, but not all problem domains succumb to these techniques. We have described a problem domain where data owners have data that they wish to share, but they cannot move to a single common ontology because of the potential loss of information (incoherency). One inadequate approach uses a series of ontologies for all possible (i.e., brut-force) combinations of data, but is prohibitively expensive.

Using our approach we have demonstrated the use of software agents to dynamically create merged ontologies, which significantly reduce the cost of developing brut-force ontologies, while providing the broadest access to distributed information. These ontologies meet the requirements stated by Gruber and others producing a shared ontology across the participating data sets. In this approach, ontologies are dynamically generated based on a user request that is evaluated against the latest ontological information from each local data source. Participating agents can each understand and provide information about all the data concepts that are shared across the participating systems, which significantly increases the capability of current (i.e., latest agent derived) ontologies. The implementation of this approach produced significant financial benefit, and will see broad deployment in the near future.

One idea that kept occurring as each new data set was added was the automation of ontology production. The idea would be that software would make an initial draft of a new data set's ontology. The time required to generate a final draft of the ontology would then be reduced, requiring only a final editing and/or verification by the data owner. We envision building ontologies by using clustering techniques; new data concepts would naturally cluster close to matching data concepts that have already been incorporated.

References

- [1] Chandrasekaran, B., Josephson, J.R. and Benjamin V.R., What are Ontologies, and Why Do We Need Them?, IEEE Intelligent Systems, Volume 14, No.1; January/February 1999, 20-26
- [2] Gruber, T., A Translation Approach to Portable Ontology Specifications, Knowledge Acquisition, 5(2), 199-220 1993
- [3] Extensible Markup Language (XML) 1.0 (Second Edition), W3C Recommendation 6 October 2000
- [4] Potok, T.E., Elmore, M.T., Ivezic N., Collaborative Management Environment, Proceedings of the InForum'99 Conference