

Project Description

The science and engineering communities are producing very large data sets that are also increasingly complex and diverse. These data sets are very well suited for particular narrowly-defined, discipline-specific purposes. In principle, these data sets could be used for solving more broadly-defined scientific problems such as understanding whole organisms, ecosystems and human populations. However, incorporating multiple data types from multiple sources to solve these problems remains a significant challenge. For example, a testable macroscopic biological hypothesis might involve the effect of environmental or climatic change on the genomic makeup of a given organism. As another example, a macroeconomic hypothesis concerning the most efficient use of resources to improve the quality of life in a region will depend on cultural and environmental knowledge as well as economic statistics.

While the data sets that are currently being developed typically engender the greatest level of enthusiasm by the communities that are creating them, data sets created in the past can have equal importance for related communities. Biodiversity is a case in point. The painstaking observations by generations of biologists over centuries represent an important resource for modern ecology and biodiversity studies, but those observations are locked in old textbooks and monographs that are not easily accessed by modern computing technology. The problem is not just the differences in recording media (paper versus disks) but also the enormous changes in terminology over time. Current data sets run the risk of an even more rapid obsolescence as the meaning of the data fields is forgotten even by the individuals who introduced them.

1 Vision and Rationale

Semantic technologies based on logic, databases and the Semantic Web can address the problem of meaningful access to and integration of data both today and over decades and centuries. This proposal is to develop and deploy a new interoperability infrastructure called the Open Ontology Repository (**OOR**) as the basis for semantic technologies. This infrastructure will serve a vibrant community of scientific researchers with collections of controlled vocabularies and knowledge models written in RDF, OWL, XML and other knowledge representation languages.

1.1 Opportunity and Impact

The OOR has the potential for having a major impact on data set interoperability in general, but to ensure that this potential is fully realized, the project will include a vigorous education and outreach program. This program will seek to bring other data-intensive research communities into the OOR initiative. Since the OOR will be an open, federated architecture and infrastructure, it is intended to be utilized by communities to host their own ontologies as well as to allow the communities to adapt previously established ontologies for their own purposes.

Yet another opportunity for a broader impact is to use ontologies as a tool for education at all levels, even at the elementary school level. Since ontologies formalize the language of a community, they can be the basis for education and training for their discipline, provided that the ontologies are properly annotated. Kenneth Baclawski is currently engaged in an NIH-sponsored research project to develop ontology-based automated health counseling tools. These tools use “relational agents” that converse with patients. Relational agents have also been developed for special-purpose educational tools. Ontologies can offer the potential for scaling up these special-purpose educational tools to apply to virtually any domain, provided that the ontology is properly annotated. The research problem is to determine what annotations are required and to minimize the cost of providing the annotations, and then to build educational tools that are based on these annotated ontologies.

1.2 Sustainable Economic and Technology Models

To address the issue of long-term sustainability, we propose to develop a new paradigm for maintaining semantic linkages available through the Internet. Specifically, we will develop a federated knowledge

repository that can collectively correct for multiple points of failure and can foster collaborative stewardship of scientific knowledge. Particular emphasis will be given to the development of technological solutions that build on existing, proven architectures for maintaining biological (e.g., BioPortal, OBO Foundry and the International Nucleotide Sequence Data Consortium) and abiotic data (e.g., the National Climatic Data Center), as well as standards for metadata and services (e.g., ISO XMDR, WSDL and UDDI).

In disciplines where opinions and facts can change over time, there is a need to retain versioning information. In biodiversity, for example, the concept of a species may be adjusted as more information (e.g., molecular) about a particular species is attained. Taxonomic hierarchies may thus change over time; however, an investigator may wish to still compile what assertions were made with respect to an “older” species concept. These types of retrospective studies can be highly informative, especially in the context of understanding how biodiversity knowledge changes over time.

The OOR and the individual repositories of the communities being served will be developed in a series of phases which will include a series of milestones to assess progress toward achieving both the technical and sustainability goals.

- **Year 1:** Gathering of requirements from the initial communities being targeted by the OOR. Development of APIs and planning for adaptation and integration of the existing ontology repository platforms. Frontier research start up in areas defined by vision. Education and outreach programs are limited to selected subcommunities.
- **Year 2:** Completion of the first release of the OOR platform. Initial delivery of OOR services to the targeted communities. Feedback from the communities will be gathered for use in the next phase. Planning begins for the second release. Frontier research focus on DataNet infrastructure enablers. Education and outreach programs expand but remain within the targeted communities.
- **Year 3:** Development of the second release begins, based on lessons learned. Education and outreach programs expand beyond the targeted communities to other DataNet partners. Delivery of OOR services to additional research communities. Frontier research focus on community interaction enablers. Framework for incorporating research results designed.
- **Year 4:** Second release of the OOR platform. Education and outreach expands to include commercial and industrial communities. Transition to self-sustaining OOR begins with planning for the organization structure. Frontier research expansion encouraged. Framework for incorporating research results integrated into the OOR.
- **Year 5:** Transition toward fully self-sustaining mode. Frontier research expansion continues and framework for incorporating results back into OOR exercised. Education and outreach are now the primary role played by the project.

The OOR team currently meets weekly using a virtual collaborative environment. This environment is managed by Peter Yim. See Section 2.4. In addition to the weekly meetings, subteams will have their own meetings, on a weekly or biweekly schedule.

2 Activities

We propose to develop an open ontology repository (OOR) of controlled vocabularies and knowledge models that have been encoded in RDF, OWL, and other knowledge representation languages. More specifically, we propose to develop an open repository for the metadata and data sets of the following communities:

- Biomedicine, including ontologies for genomics, proteomics, diseases, anatomy, model organisms, and other resources served by the highly successful BioPortal repository.
- Biodiversity, especially the evolutionary history and taxonomic communities that strive to create syntheses of information such as the species pages in the Encyclopedia of Life.
- Climate and environmental communities (including both natural environments and built environments).

- Human culture and society, with a focus on ensuring compliance with federal regulations for privacy and security across diverse data collections.

2.1 Support the full data preservation and access lifecycle

To truly support discovery, innovation and learning well into the future the OOR collaboration will manage the full data life cycle by providing an architecture and an infrastructure that supports a) the creation, sharing, searching, and management of ontologies, and b) linkage to database and XML Schema structured data and documents. Complementary goals include fostering the ontology community, the identification and promotion of best practices, and the provision of services relevant to ontologies and instance stores. Examples of anticipated services include automated semantic interpretation of content expressed in knowledge representation languages, the creation and maintenance of mappings among disparate ontologies and content, and inference over this content. The OOR will support a broad range of semantic services and applications of interest to enterprises and communities.

The OOR will develop efficient logic programming-based reasoning methods that amalgamate Semantic Web-based ontologies and rules with extended Prolog and Answer Set Programming, to be used for reasoning over the ontologies, instances, and rules of the repository. [57, 56, 68] The OOR will design and implement service-oriented architectures and services, including automated and semi-automated service orchestration and parallel optimization to support the repository. [40, 53, 69]

The following are the core requirements for the OOR:

1. The repository architecture shall be scalable.
2. The architecture shall be optimized for sharing, collaboration and reuse.
3. The repository shall be capable of supporting ontologies in multiple formats and levels of formalism.
4. The repository architecture shall support distributed repositories.
5. The repository architecture shall support explicit machine usable/accessible formal semantics for the meta-model of the repository.
6. The repository shall provide a mechanism to address intellectual property and related legal issues/problems.
7. The repository architecture shall include a core set of services, such as support for adding, searching and mapping across ontologies and data related to the stored ontologies.
8. The repository architecture shall support additional services both directly within the province of the repository and as external services.
9. The repository should support all phases of the ontology lifecycle.

2.1.1 Data deposition/acquisition/ingest

The OOR will develop development of requisite ontology-based architectures, including ontology lifecycle management, theories and implementations of ontology modularity, upper and middle ontologies, and research and software development of methods for automatically and semi-automatically aligning and mapping ontologies. Logical relationships between ontologies will be supported within the repository, including mutual consistency, extension, entailment, semantic mappings, intelligent search, and decision support. [13, 42, 47, 64, 66, 65]

The OOR will support internal and external services and applications including: ontology creation tools, ontology editors, ontology differencing tools, ontology modularization tools (clustering, etc.), ontology export, ontology visualization (e.g., graph visualization), version management and access control. While the emphasis is on the metadata level, ontologies also include instance data. For knowledge-rich domains such as the targeted sectors, the ontology includes all of the data as well as the metadata. For other domains, the data will be managed by special-purpose applications, and the ontology will be only part of the database, playing the roles that are most appropriate, such as encoding access policies and procedures, enabling discovery and interoperability, and ensuring that data remains accessible and understandable over timelines of decades or more.

2.1.2 Data curation and metadata management

We distinguish between gatekeeping and quality control. Gatekeeping criteria are a set of minimal requirements that any ontology within the OOR has to meet. The latter are intended to enable the users of the OOR to quickly find ontologies that fit their needs; the criteria are not supposed to ensure the quality of the ontologies. The ontologies in the OOR must meet the following criteria: (1) The ontology is submitted in a publicly described language and format; (2) The ontology is read accessible; (3) The ontology is expressed in a formal language with a well-defined syntax; (4) The authors of the ontology provide the required metadata; (5) The ontology has a clearly specified and clearly delineated scope; (6) Successive versions of the ontology are clearly identified; (7) The ontology is appropriately named; It is especially important that the required metadata include information about the process that is employed to create and maintain the ontology. (Is the ontology maintained in a cooperative and transparent process? Can anybody participate in this process?) Further, the metadata has to include information about the license under which the ontology is submitted.

In addition to the logical metadata for ontologies, the OOR will include metadata for ontologies considered as engineering artifacts. This includes provenance, versioning, existing applications of the ontology (e.g. interoperability, search, decision support) and domain-specificity (e.g. biology, supply chain management, manufacturing).

The Ontology Metadata Vocabulary (OMV), Dublin Core, ISO 11179, ISO 19763, and other existing approaches to provenance and versioning metadata will be used as the basis for the metadata for ontologies in the OOR. An empirical approach will be used to identify and evaluate ontology metadata. Proposals for ontology metadata already exist, and they will be evaluated using use-case scenarios. These scenarios both motivate the use of the metadata and help establish best practices.

2.1.3 Data protection

An ontology repository requires mechanisms for effective management. The understanding is that as a repository and its infrastructure evolve, more management support mechanisms will be included. The core mechanisms to be provided in the first version include enforcement of policies for access, submission, governance, change management, and control over user and administrative access. The later version of the OOR will provide capabilities to: create usage reports, validate syntax, check logical consistency and automatically categorize a submission.

2.1.4 Data discovery, access, use, and dissemination

Achieving these goals will help reduce semantic ambiguity whenever and wherever information is shared, thereby allowing information to be located, searched, categorized, and exchanged with a more precise expression of its content and meaning. The artifacts of the repository will provide a semantic grounding for diverse formats and domains, ranging from the conceptual domains and specific disciplines of communities to technical schemas such as WSDL, UDDI, RSS, and XML schemas, and of course expressed in standard ontology languages such as RDF, OWL, Common Logic, and others. Perhaps most importantly, the repository will enable wide-scale knowledge re-use and reduce the need to re-invent the wheel when defining concepts and relationships that are already understood.

To facilitate knowledge discovery the repository shall provide metadata capabilities to support search capabilities, governance process, and management. The repository will support discovery by at least the following criteria: domain, author/creator/source, version, language, terminology and controlled vocabularies, quality, mapping, and inference. The interfaces for discovery will be suitable for both specialist and non-specialist users, using GUIs, web services, and language-based APIs.

2.1.5 Data interoperability, standards, and integration

To support the sharing and reuse of ontologies within the repository the OOR will store both ontologies and metadata for ontologies. The metadata will allow users to: (1) determine whether an ontology is suitable

for a user purpose; (2) capture the design rationales that underlie the ontology; (3) find information about author, author credentials, and source of ontology reference material; (4) retrieve ontologies for use in domain applications; (5) retrieve ontologies to be integrated with other ontologies; (6) retrieve ontologies that will be extended to create new ontologies; (7) determine whether or not an ontology can be integrated with given ontologies; (8) determine whether a set of ontologies retrieved from the repository can be used together; and (9) determine whether an ontology in the repository can be partially shared.

There will be policies for creation and modification of metadata and documentation of ontologies and the management of the persistence and sustainability of ontologies.

Users (including end-users, ontology and repository developers, subject matter experts and stakeholders) should participate in the collaborative ontology development life cycle and in decisions regarding what metadata are suitable for ontologies in the repository.

The metadata will include both logical metadata (logical properties of the ontology independent of any implementation or engineering artifact) and engineering metadata (properties of the ontology considered as an engineering artifact).

2.1.6 Data evaluation, analysis, and visualization

It is not sufficient for the OOR just to store ontologies, it also needs to enable the evaluation of the ontologies within it. The OOR will offer functionalities like those on social networking sites which would allow users to comment on ontologies and rank them. Further, the OOR will enable selective views of the repository using tags provided by subcommunities that characterize ontologies with respect to their chosen criteria. For example, such a view might select ontologies for specific fields of research or industries, or for ontologies satisfying specific quality criteria or levels of organizational approval.

The OOR will develop methods, practices, services, and artifacts to support automated and human reviewed evaluation and comparison of ontologies stored in the repository. [34, 45, 44]

2.2 Engage at the frontiers of science and engineering research and education

Ontology engineering is a rapidly developing research area. There are both computational and storage challenges as the size of a knowledge base grows, since the computational complexity of logical inference is much greater than the complexity of traditional database query processing. The OOR team is closely connected both with the ontology engineering community and with the communities specifically targeted by this proposal. As shown in the biosketches, the team members have substantial ontology research experience. Ontology engineering research will be an integral part of the OOR effort.

2.3 Education and training

As an integral part of the proposed project, the OOR will support a vigorous educational outreach program to bring other data-intensive research communities into the OOR initiative. Since the OOR will be an open, federated architecture and infrastructure, it is intended to be utilized by communities to host their own ontologies as well as allowing the communities to adapt previously established ontologies for their own purposes.

The Ontolog Forum has been engaging in educational and outreach activities for 6 years, reaching over 40 distinct communities. Examples include communities in bioinformatics, national command and control, and intelligence. [36, 35, 33, 48, 46, 51, 64, 67, 68, 69]

2.4 International community and user input and assessment

This OOR initiative has actually emerged from members of a community that has successfully developed an infrastructure which has grown and sustained itself since early 2002. Ontolog (or the Ontolog Forum) is an open virtual community for ontology engineering, with members who realize the importance and potential impact of ontology, and are passionate about moving it into the mainstream through adoption and

standardization. With over 550 participating members from around 30 different countries, Ontolog has been able to amass a highly regarded body of knowledge in the ontology domain. This was done by archiving (and metadata tagging) all its membership dialog and contributions transacted over its collaborative work environment infrastructure, which also serves as a dynamic knowledge repository for the community's collective intelligence. This open repository receives about 50,000 page views or file downloads per day, from around 120 cities worldwide [70].

The "Open Ontology Repository" theme was selected by the community as the theme for their 2008 Ontology Summit. This is an annual event (since 2006) jointly organized by NIST, Ontolog and the National Center for Ontological Research. The Ontology Summit entailed a three-month online discourse, four virtual panel discussions, a two-day face-to-face workshop at NIST, and the publication of a Communique. The initiative was co-sponsored by about 50 institutions, and had an Organizing Committee and an Advisory Committee made up of more than 40 individuals, who are among the most respected names, worldwide, in the ontology and ontological engineering domain [49].

3 Organizational Structure

3.1 Leadership and Management

The OOR initiative intends to designate Peter Yim as Director of the project. Peter brings along an illustrious visionary, strategic leadership and cross-disciplinary management portfolio. He has managed as many as 3000 people and 300 software engineers, and has built companies from scratch to \$500 million in revenue. He has learned the art and science of working effectively at almost any part of a real or virtual organization – from the executive office, to corporate boardrooms, to R&D, all the way to the shop floor. He works in corporate settings and academia, as well as on non-profits, education and government boards. He has been a co-founder of Ontolog Forum (along with Leo Obrst); the OASIS Universal Business Language (UBL) standard technical committee; Director of Program Management for VerticalNet, Inc., the first B2B company that employed ontologies for data integration in an eCommerce setting; as well as having been a principal or chief executive at other industrial enterprises. He now heads CIM3, a distributed collaboration technology and internet service company which has been providing infrastructure support to the Ontolog Forum operations, as well as to various US Federal inter-governmental collaborative endeavors.

The OOR project shall be structured into sub-teams with appropriate sub-team leadership to handle the R&D in areas like infrastructure, integration, content, tools and community. The sub-team leadership reports to a leadership team headed by the project director. The OOR initiative will also establish an Advisory Board where prominent figures in the relevant technology domains will be invited to join and advise the team.

3.2 Comprehensive expertise and infrastructure capacity/capabilities

The OOR team possesses considerable expertise in library and archival sciences (Mark Musen, Indra Neil Sarkar); computer, computational and information sciences (all team members); cyberinfrastructure (all team members); and domain sciences (discussed in Section 3.3 below). The team has been collaborating at virtual meetings for over a year, as well as at a face-to-face meeting at NIST in April of 2008. The team has developed an effective organizational structure that enables shared responsibility, close coordination and cooperation, and catalyzes the rapid exchange of ideas. The team has access to considerable computational, storage, network access, dissemination, interaction and communication resources. Peter Yim is responsible for managing this infrastructure as explained in the previous section above.

3.3 Diverse, multi-sector participation

The initial sectors that will be served by the OOR include the following: (1) **Biology**, (2) **Biodiversity**, (3) **Climate and Environment**, (4) **Human culture and sociology**.

BioPortal is a centralized repository for biomedical ontologies. It primarily serves the biomedical research sector, including the genomics, proteomics, diseases, anatomy, and model organism communities. However, BioPortal is a general purpose ontology repository, and it will be the foundation for the OOR. Mark Musen and his team at the NCBO will provide the expertise for this sector.

Indra Neil Sarkar and his team at MBL will provide the expertise for serving the biodiversity community. The focus will be on the evolutionary history and taxonomic communities that strive to create dynamic syntheses of information such as the species pages in the Encyclopedia of Life. He will also be providing expertise in the natural environment sector.

The Marine Metadata Interoperability (MMI) project team will collaborate with the OOR project to integrate MMI project work with the OOR effort. Coordinated by Principal Investigator John Graybeal of the Monterey Bay Aquarium Research Institute, MMI has as its key mission objective developing a broad community presence to address marine metadata issues. The MMI project is uniquely positioned to understand community needs, address those needs in open and broadly applicable ways, and obtain community engagement and buy-in for open solutions. With Technical Lead Carlos Rueda and contributor Luis Bermudez (Southeastern Universities Research Association), MMI's pioneering work developing a community semantic architecture and ontology repository will provide key insights and building blocks for the OOR project. MMI's marine science solutions are equally valuable and viable in most environmental communities, and a natural progression of the effort; several environmental science communities are represented at our semantic interoperability workshops.

The climate and man-made environment is a diverse sector that includes architecture and engineering communities. Katherine Goodier and her colleague Thomas Lyndon Wheeler will be providing expertise for this sector as well as the human culture and sociology sector.

For each of these sectors, the OOR collaboration will be providing: (1) Guidance for data providers, metadata providers, and ontology providers; (2) Organized references on all facets of metadata needs and solutions; (3) Services targeting semantic interoperability in the respective and related domains, including vocabulary lists, ontology repository and associated services, and vocabulary creation and maintenance tools, services, and guidance; (4) Community collaboration environment (shared files, email archives, and web pages, either public or secure); (5) Access to work in progress on metadata tasks and projects. In addition, the OOR collaboration can provide purposeful capabilities: (1) Targeted identification and evaluation of resources (vocabularies, standards, tools, services); (2) Identification and engagement of key community participants (projects or individuals) in metadata initiatives; (3) Training and workshops in metadata technologies and techniques, particularly dealing with semantic tools and services, including vocabulary and ontology development, metadata standards and their application, as well as metadata-enlightened architectural development and analysis; (4) Community environment(s) to advance particular topics or discussions. The references section of the proposal includes links to resources that are already available in the targeted sectors.

3.4 Data Network

The primary purpose of ontologies is to achieve interoperability. The OOR will initially be focused on collaboration and coordination with its target communities. As the project continues, it will collaborate and coordinate closely with other DataNet Partners and digital preservation/access organizations both nationally and internationally. The ultimate goal is to allow for seamless, single entry point discovery, access, and use of data from any source. The Ontolog Forum in general, and the OOR team in particular, are already engaged in developing and disseminating best practices and principles. The Ontolog Forum has a substantial, multi-year track record for shared governance and coordination that will be leveraged for the proposed project.