# Harnessing the Power of the Community in a Library of Biomedical Ontologies

Natalya F. Noy, Michael Dorf, Nicholas Griffith, Csongor Nyulas, Mark A. Musen

Stanford University, Stanford, CA 94305, US
noy,mdorf,ngriff,cnyulas,musen@stanford.edu

**Abstract.** Biomedical ontologies provide essential domain knowledge to drive data integration, information retrieval, data annotation, natural-language processing, and decision support. The National Center for Biomedical Ontology is developing BioPortal, a web-based library of biomedical ontologies. As the biomedical community develops the ever growing set of ontologies of varying sizes, quality, purpose, level of logical rigor, it becomes more and more difficult to find the "best" and "most appropriate" ontologies in the domain. In BioPortal, we use the social approaches in the Web 2.0 style to bring structure and order to the collection of biomedical ontologies. BioPortal enables users to provide and discuss a wide array of knowledge components, from submitting the ontologies themselves, to commenting on and discussing classes in the ontologies, to reviewing ontologies in the context of their own ontology-based projects, to creating mappings between overlapping ontologies and discussing and critiquing the mappings. In this paper, we discuss the community features of the BioPortal ontology library and describe the infrastructure that supports these features. BioPortal is available online at http://bioportal.bioontology.org.

## 1 The library of biomedical ontologies in BioPortal

As the number of ontologies available for Semantic Web applications grows, so does the number of ontology libraries that index and organize the ontologies. Some libraries get the ontologies by crawling the Web (e.g., Swoogle [4], Watson [3] and OntoSelect [1]). In other libraries, users submit their ontologies themselves (e.g., the Protégé ontology library[1]). Some libraries provide strict selection criteria for inclusion (e.g., OBO Foundry [14]). All these libraries provide a gateway for users and application developers who need to find ontologies to use in their work. In our laboratory, as part of the National Center for Biomedical Ontology (NCBO), we have developed BioPortal[2]—an open library of biomedical ontologies. Researchers in biomedical informatics submit their ontologies to BioPortal and community members can access the ontologies in their web browsers through the BioPortal user interface or through web services [12]. The BioPortal users can browse and search the ontologies, update the ontologies that they authored by uploading new versions, comment on any ontology (or portion of an ontology) in the library, evaluate it, describe their experience in using the ontology, or make suggestions to ontology developers. This focus on enabling members of the

---

[1] http://protegewiki.stanford.edu/index.php/Protege_Ontology_
Library
[2] http://bioportal.bioontology.org

community to contribute actively to BioPortal content and to increase its value to other users, distinguishes BioPortal from other ontology libraries.

At the time of this writing, BioPortal has 160 ontologies, covering a wide range of domains in biomedicine, from anatomy, to diseases, to protein descriptions. BioPortal supports several formats for representing ontologies: the OBO format, Rich Release Format (RRF) from the US National Library of Medicine (for terminologies comprising the Unified Medical Language System [8]), OWL, RDF(S), and the Protégé frames format. BioPortal uses the Mayo Clinic's LexGrid system[3] to store ontologies in OBO Format and to access standard biomedical terminologies, such as UMLS in RRF. Protégé[4] serves as the backend for OWL and RDF ontologies.

This paper makes the following contributions:

– We define several types of contributions that a community-based ontology library can collect and aggregate, including ontologies and their successive versions, notes and threaded discussions on classes, reviews of ontologies in the context of specific projects, ontology mappings, and views and subsets of ontologies.
– We describe the validation of our approach in the form of a production implementation of BioPortal.

## 2  Community Features of BioPortal

With the open, community-based nature of the BioPortal library of biomedical ontologies, we are experimenting with the new ways of publishing, evaluating, and integrating the knowledge infrastructure that is essential to life sciences today. Specifically, BioPortal users and ontology developers can contribute a variety of information to the library, increasing its value to others. First, the library enables all members of the biomedical community to **publish** their ontologies (Section 2.1). Second, other members of the user community can provide **feedback** on specific elements of the ontologies or ontologies as a whole (Section 2.2). Third, ontology users can describe their experiences in using the ontologies from the library in their ontology projects, thus providing a novel way of **evaluating** ontologies (Section 2.3). Fourth, we open the process of declaring relationships between concepts in different ontologies (**ontology mapping**) to the community by enabling users to declare mappings between concepts and to comment on mappings created by others (Section 2.4). Finally, users can publish and describe **subsets** or different **views** of the ontologies in BioPortal (Section 2.5).

### 2.1  Publishing Ontologies in BioPortal

Any developer of an ontology that is relevant to biomedical domain can publish it in BioPortal. When submitting an ontology to BioPortal, the user must provide essential metadata about the ontology, such as its name and acronym, the domain that the ontology covers, keywords, links to additional information, and provenance information, including who developed the ontology, version details, dates of the release, and so on.

---

[3] http://informatics.mayo.edu/LexGrid
[4] http://protege.stanford.edu

The ontology authors then have two choices in terms of submitting the ontology itself. They can choose not to submit the ontology content, thus providing BioPortal only with the ontology metadata. BioPortal users will then be able to see the ontology metadata in the library and to comment on the ontology as a whole. They will not be able to view and search the ontology content though. We designed this submission option for ontology authors who are not willing to make their ontology accessible to the community directly through BioPortal (e.g., because of licensing issues) but would still like the community to know about their resource. Most ontology authors, however, choose the second option: submit both the metadata and the ontology itself. After the ontology author submits an ontology, BioPortal parses and indexes it and makes it available for searching and browsing. All ontologies in BioPortal are publicly accessible.

**Ontology Web Services and Links.** When ontology authors submit their ontology to BioPortal, they enable a wider user community to find and use their ontology. Furthermore—and in some cases more important—publishing in BioPortal is a very easy way to get a web presence for an ontology. Many ontology developers are neither interested nor willing to host their own web server. By uploading their ontology to BioPortal, they get a web link (a URL) that they can give to their users when inviting these users to browse the ontology, see details of specific concepts, visualize the ontology or any of its parts. BioPortal also provides URLs to access any class in the ontology directly. Thus, for example, an ontology author can send to her collaborators a URL for a specific class that she wants to discuss with them. The collaborator can then see the details of the class definition and all the related information by following the link. Any BioPortal user can subscribe to an RSS feed of changes to a specific ontology in order to get notified of any user-contributed content relevant to that ontology, such as comments or mappings.

In addition to accessing an ontology and its components in a web browser, users can use the BioPortal RESTful API to access any ontology or its components through a web service. In fact, the BioPortal user interface itself uses this REST API to display most of the information that the users see on the BioPortal web site. There are web services to get metadata about an ontology, its root classes, details of any concepts, hierarchical information for any concept; there are web services to download an ontology, get a diff between two versions, get notes or mappings for an ontology [12].[5]

**Ontology Versioning.** Most, if not all, ontologies in BioPortal continue to evolve and authors continue to publish new versions. Thus, any ontology library must address the issue of ontology evolution [9]. Ontology authors can submit successive versions of their ontologies to BioPortal. Each version can have its own set of metadata, since any detail about an ontology (from its scope, to provenance details, to relevant links) can change from one version to another. Users can explore and use any version of any ontology in the collection (Figure 1). For each ontology, BioPortal provides two sets of services and links: one set resolves to a specific version requested by the user; another set resolves to the latest version of an ontology. Thus, if a user's application relies on a

---

[5] See `http://bioontology.org/wiki/index.php/NCBO_REST_services` for the details of REST services that NCBO currently offers.

**Fig. 1. Ontology details in BioPortal.** The details page describes the information about the ontology, its provenance, lists the versions of the ontology and provides an overview of concepts that have notes and mappings.

specific version, they can pass that version as the parameter to the services. If the user wants to get the latest version, whatever that version is, they can pass in the generic ontology id (we call it a "virtual" id); BioPortal redirects this call to use the latest version of the ontology.

**BioPortal and OBO Foundry** The OBO Foundry initiative [14] aims at creating a set of well-documented and well-defined ontologies that are designed to work with one another. The OBO Foundry has an editorial process defining which ontologies become part of that collection. For many ontologies that are "OBO Foundry candidates" the OBO Foundry site is the primary publication vehicle. Because of the importance of this collection to the biomedical community, BioPortal includes all the ontologies from OBO Foundry in its collection. We have an automatic process that checks the OBO Foundry site nightly and pulls in the updates to ontologies and ontology metadata into BioPortal. Thus, the BioPortal collection includes all the OBO Foundry ontologies. These ontologies currently constitute about 40% of the BioPortal collection.

### 2.2 Providing Comments

Users can add notes to classes in BioPortal, discussing the rationale for modeling decisions, pointing out problems with definitions, requesting changes from ontology authors, and so on. These notes are attached to specific classes and one can think of them as metadata on those classes. Notes can be organized in a threaded discussion (Figure 2).

So far, we have observed a variety of use cases for notes, including passing the feedback on the classes to the ontology authors, suggesting changes and corrections,
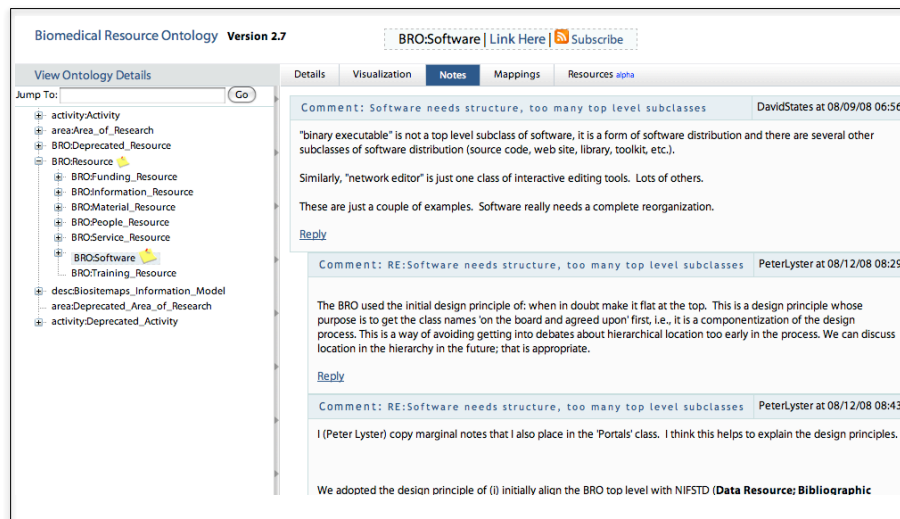
**Fig. 2.** A class hierarchy and a threaded discussion about a class in BioPortal.

requesting new items, discussing class definitions among a group of contributors, and providing additional information about a class, such as references, images, or supporting documentation.

### 2.3 Ontology Evaluation as a Community Process

One of the main functions of an ontology library in general and of BioPortal in particular is to help researchers find the ontologies that they can use in their applications. When confronted with a list of published ontologies in his domain of interest, the key question for the user in helping him decide which ontology to use is the following: "Has anyone used this ontology successfully for a task that is similar to mine?" In the BioPortal project, we are experimenting with the social ways of evaluating ontologies and answering this question.

Researchers have proposed a number of ways to evaluate an ontology (e.g., [6, 5]). These approaches evaluate the structure of the ontology and its consistency or conformance to certain principles. While these approaches can be very helpful in understanding whether an ontology is structured properly, they do not provide much insight on how well an ontology covers a particular domain, or how appropriate it might be in certain types of applications. Thus, we have developed an infrastructure that would enable ontology *users* to provide this additional—and often essential—information.

BioPortal enables its users to describe their ontology-based **projects**. After a user describes a project, he can select the BioPortal ontologies that his project uses and then provide the reviews of these ontologies *in the context* of his project. For example, an ontology that is an excellent resource, say, for an information-extraction application

because the ontology contains lots of lexical information, may not be appropriate for an application that needs to use an ontology for reasoning if its classes do not have axioms in their definitions.

Each ontology review has different dimensions. We have conducted surveys of BioPortal users to determine which review dimensions they would find particularly useful. We currently have the following dimensions for reviews: domain coverage; correctness; quality of content; degree of formality; documentation and support; usability. Each review refers to a specific version of an ontology and includes a star rating and a text description for each of the dimensions.

## 2.4 Community-Based Ontology Mapping

Ontologies in BioPortal, as in almost any ontology repository, overlap in coverage. Thus, **mappings** among ontologies constitute a key component that enables the use of the ontologies for data and information integration. For example, researchers can use the mappings to relate their data, which had been annotated with concepts from one ontology, to concepts in another ontology. We view ontology mappings as an essential part of the BioPortal library. In BioPortal, users can browse the mappings, create new mappings, upload the mappings created with other tools, download mappings that BioPortal has, or comment on the mappings and discuss them [10].

Our implementation enables and encourages *community participation in mapping creation*. We enable users to add as many or as few mappings as they like or feel qualified to do. Users can use the discussion facilities that we integrated in BioPortal to reach consensus on controversial mappings or to understand the differences between their points of view. Most researchers agree that, even though there has been steady progress in the performance of the automatic alignment tools [2], experts will need to be involved in the mapping task for the foreseeable future. By enabling community participation in mapping creation, we hope to have more people contributing mappings and, hence, to get closer to the critical mass of users that we need to create and verify the mappings. The BioPortal mapping repository contains mostly the mappings created by our users elsewhere and by other tools, and uploaded in bulk to BioPortal.

With this large number of mappings coming from different sources, we expect that different users and algorithms would map one concept from an ontology $O_1$ to different concepts in another ontology $O_2$. Our infrastructure supports this plurality of mappings and we plan to use social means to determine the "best" mappings or the mappings that would be more appropriate in one context and may not be appropriate in another.

In a repository where users can contribute data, enabling extensive metadata for mappings is critical. Thus, for each mappings we store the metadata on what the source of the mapping is, how the mapping was created and in which application context, which algorithm, if any, was used to create the mapping, and which version of the algorithm and which configuration parameters, who uploaded it to BioPortal and when.

At the time of this writing, the BioPortal mapping repository contains more than one million mappings, the majority of which were created using various automatic or semi-automatic algorithms and then uploaded to BioPortal.

### 2.5 Views, Subsets, and Value Sets

Finally, users can upload views or subsets of BioPortal ontologies. A view can be a subset of ontology concepts that was created for a particular purpose. For instance, BioPortal contains NeuroFMA, a subset of FMA classes relevant for neuroimaging. A rendering of an ontology in another format can also be represented as a view. Indeed, if a user translates one of the BioPortal ontologies into a different language (e.g., Chinese), that user can upload the translated ontology as a view on the original one.

While ontology authors or administrators who upload the ontology to BioPortal control which new versions get uploaded and when, anyone in the community can contribute **views** on any ontology in BioPortal. A view represents a materialized subset of an ontology created for a particular purpose [13]. One can consider a view to be just another ontology in BioPortal that has additional metadata describing which "master" ontology (and which version of it) was used to create the view, how the view was created (e.g., the specific query and engine that was used to extract it), what was the purpose for creating the view, and so on. Because we represent views simply as ontologies in BioPortal (albeit with special status), users can review the views, comment on their use in their projects, discuss where a particular view is appropriate and so on.

## 3   BioPortal Implementation of Community-Based Features

We describe the details of our internal representation elsewhere [13] and we present it here briefly to describe what happens "under the hood" when BioPortal supports the community-based features that we described in Section 2.

We use an ontology-based approach to represent all the metadata in BioPortal. In this context, we refer to all the data *about* the ontologies (e.g., ontology details, comments, reviews, mappings) as metadata. We use an ontology—the BioPortal Metadata Ontology—to describe the structure of the metadata and the metadata values themselves are represented as instances in this ontology.

The BioPortal Metadata Ontology is an OWL ontology that imports a number of other ontologies (Figure 3) and includes classes to describe an ontology itself, its versions, metadata properties about the ontology, creators of an ontology, user-contributed content, such as notes, reviews, mappings, and views.

The BioPortal Metadata Ontology imports several ontologies that deal with the types of metadata that BioPortal supports:

- **The Ontology Metadata Vocabulary (OMV)** describes most of the metadata for ontologies themselves (e.g., domain, author, version number, ontology language, etc.). The OMV provides the vocabulary for describing a specific ontology version. An instance of the class `OMV:Ontology` describes a single version of an ontology. This class contains properties describing pertinent information about the ontology in general.
- **The Protégé Changes and Annotations Ontology** (CHAO) provides the definitions for generic annotations (the `Annotation` class) and ontology components that they annotate. We use the instances of the Protégé CHAO ontology to represent comments that BioPortal users contribute to the ontologies. Each comment
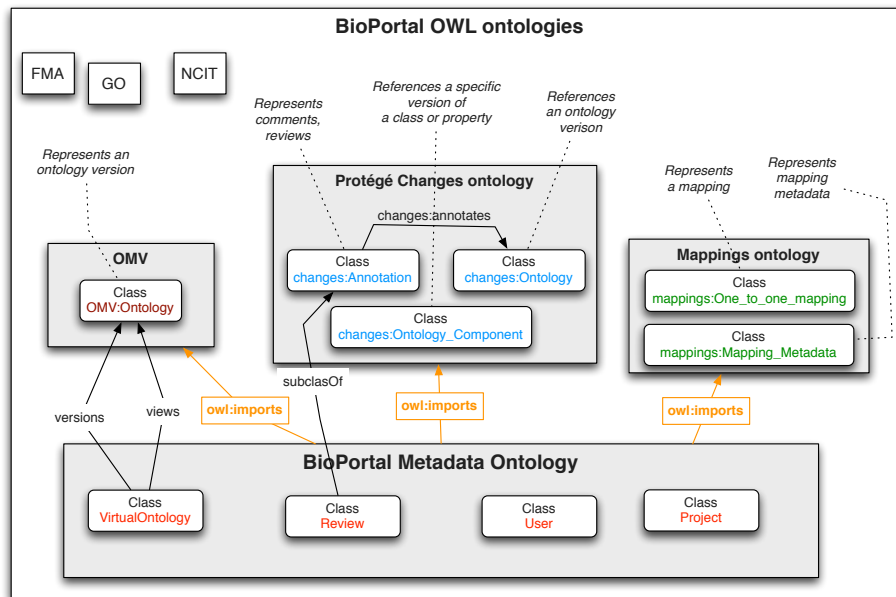
**Fig. 3. The BioPortal Metadata Ontology**: Some key classes and ontologies that the BioPortal Metadata Ontology imports. The BioPortal Metadata Ontology is itself it the BioPortal ontology repository, along with domain ontologies such as the Gene Ontology (GO), the FMA, and others.

is represented as an annotation attached to a specific class (in a specific ontology version) or to another annotation (if it is a response to a comment). The same mechanism exists in Collaborative Protégé, a version of the Protégé ontology editor that supports collaborative ontology editing [15]. Because BioPortal and Protégé share the same structure for representing user comments and discussions, one can potentially open a BioPortal ontology in Protégé and see the comments contributed by BioPortal users.

– **The Protégé Mapping Ontology** provides vocabulary for describing one-to-one mappings between concepts and the mapping metadata.

## 4   Discussion and Challenges

We have described a set of community-based features that we have implemented in BioPortal, a library of biomedical ontologies. In implementing these features, we are defining new models for publishing, evaluating, and integrating biomedical ontologies. And while our implementation provides the infrastructure to validate the efficacy of these ideas, we haven't had the tools released for sufficient period of time yet to validate the social component of our experiment. The main question remains "if we build it, will they come?" Our initial indicators are encouraging, with the size of the BioPortal

repository doubling in the last few months, large sets of mappings being submitted by several users, and the number of BioPortal users steadily growing. However, the community-based approaches become much more valuable when the community grows in size. Over the coming months, we plan to evaluate which features are more popular with our users and to improve our support for those features.

Our work on community-based approaches to publishing and maintaining biomedical ontologies also highlights several research issues and challenges.

As we noted earlier, ontologies inevitably *evolve* and authors publish new versions. Thus, we must maintain all the metadata, notes, reviews, and mappings through this evolution process. Users add the metadata for specific ontology versions, and, in theory, any metadata can get invalidated when a new version is published. For instance, if a class definition changes, a mapping may become invalid; or a note, requesting a change to a class, is no longer relevant. Similarly, a review that indicates some problems with an ontology may no longer be relevant after the ontologies has been fixed. At the same time, we do not want to invalidate all the user-contributed content linked to an ontology once a new version of that ontology is uploaded: Our earlier research shows that only a small fraction—usually 1-4%—of ontologies changes from one version to the next [11]; thus, a large portion of the user-contributed content is relevant for the new version. Our current approach to maintaining metadata through ontology evolution is a hybrid one: all metadata, such as comments, mappings, reviews, are attached to a specific ontology version. However, the metadata also references the global ("virtual") ontology id and the user interface exposes the metadata when users access a newer ontology version. In the future, we plan to add a subtle cue that indicates that the metadata item was created for an earlier version (since we already have that information). Furthermore, we plan to add mechanisms for archiving metadata that may no longer be relevant. For the archiving of metadata, however, we must develop policies on who has the right to archive a comment or a mapping that is no longer relevant. The authors of that comment, mapping, or review? The authors of the ontology? Only the BioPortal administrators? There are good arguments for the validity of any of these choices and we plan to discuss with our user community which approach would be the most meaningful in our case.

As with any initiative that is open to contributions from a wide variety of users, *trust* is a critical issue that we must address. While we do not have the problem of having to filter out the content from malicious or incompetent users at the moment, this problem will inevitably arise if BioPortal is successful. We plan to use an open-rating system and a web of trust [7] to enable users to rate not only the content of the repository but also the ratings and reviews (similar to, say, reviews on Amazon). Thus, we will build a web of trust network among our users. Some of the opinions about ontologies are subjective, in part because workers use ontologies for different purposes and thus value different types of feature. Therefore, we envision that users will tend to select other users with similar interests and requirements in their web of trust.

We also learned that having all information in BioPortal available to all users may be a problem for some communities. For instance, there can be a community of users that wants to discuss their ontology, or maybe test some mappings in a "private" space before making it publicly available. Thus, we are working on implementing group-

specific views of some BioPortal content, enabling groups to have discussions in private before publishing their results to the broader research community.

Finally, one of the most challenging, but also most interesting issues, is *evaluating* the contribution of our work. We are currently working on the protocols that we can use to assess the effect of community-based evaluation on the process of ontology selection. We are analyzing the mappings contributed by different users to evaluate the degree of overlap between the mappings and the degree of agreement between them. We are working on developing and evaluating new ways in which structured notes can facilitate collaborative ontology development.

## 5 Acknowledgments

## References

1. P. Buitelaar, T. Eigner, and T. Declerck. OntoSelect: A dynamic ontology library with support for ontology selection. In *Demo Session at the Intl. Semantic Web Conf. (ISWC 04)*, Hiroshima, Japan, 2004.
2. C. Caracciolo, et.al. Results of the ontology alignment evaluation initiative 2008. In *3d Int. Workshop on Ontology Matching (OM-2008) at ISWC 2008*, Karlsruhe, Germany, 2008.
3. M. d'Aquin, C. Baldassarre, L. Gridinoc, S. Angeletou, M. Sabou, and E. Motta. Watson: A gateway for next generation semantic web applications. In *Poster session at the International Semantic Web Conference (ISWC 2007)*, Busan, Korea, 2007.
4. L. Ding, et. al. Swoogle: A search and metadata engine for the semantic web. In *13th ACM Conference on Information and Knowledge Management (CIKM'04)*, Washington DC, 2004.
5. A. Gangemi, C. Catenacci, M. Ciaramita, and J. Lehmann. Modelling ontology evaluation. In *Proceedings of the Third European Semantic Web Conference*. Berlin, Springer, 2006.
6. N. Guarino and C. Welty. Evaluating ontological decisions with OntoClean. *Communications of the ACM*, 45(2):61–65, 2002.
7. R. Guha, R. Kumar, P. Raghavan, and A. Tomkins. Propagation of trust and distrust. In *13th Intl. Conf. on World Wide Web (WWW-04)*, pages 403–412, New York, NY, 2004. ACM.
8. D. Lindberg, B. Humphreys, and A. McCray. The unified medical language system. *Methods of Information in Medicine*, 32(4):281, 1993.
9. N. F. Noy, A. Chugh, W. Liu, and M. A. Musen. A framework for ontology evolution in collaborative environments. In *Fifth International Semantic Web Conference, ISWC*, volume LNCS 4273, Athens, GA, 2006. Springer.
10. N. F. Noy, N. Griffith, and M. A. Musen. Collecting community-based mappings in an ontology repository. In *7th Intl. Semantic Web Conf. (ISWC 2008)*, Germany, 2008.
11. N. F. Noy and M. A. Musen. Ontology versioning in an ontology-management framework. *IEEE Intelligent Systems*, 19(4):6–13, 2004.
12. N. F. Noy, et.al. BioPortal: ontologies and integrated data resources at the click of a mouse. *Nucleic Acids Research*, 10.1093/nar/gkp440, 2009.
13. C. I. Nyulas, N. Noy, M. Dorf, N. Griffith, and M. A. Musen. Ontology-driven software: What we learned from using ontologies as infrastructure for software. Tech report BMIR-2009-1382, Stanford University, 2009.

14. B. Smith, et.al.. The OBO Foundry: coordinated evolution of ontologies to support biomedical data integration. *Nature Biotechnology*, 25(11):1251–5, 2007.
15. T. Tudorache, N. F. Noy, S. Tu, and M. A. Musen. Supporting collaborative ontology development in protege. In *7th Intl. Semantic Web Conf. (ISWC 2008)*, Germany, 2008.