

Investigating Term Reuse and Overlap in Biomedical Ontologies

Maulik R. Kamdar, Tania Tudorache, and Mark A. Musen

Stanford Center for Biomedical Informatics Research, Department of Medicine, Stanford University

ABSTRACT

We investigate the current extent of term reuse and overlap among biomedical ontologies. We use the corpus of biomedical ontologies stored in the BioPortal repository, and analyze three types of reuse constructs: (a) explicit term reuse, (b) *xref* reuse, and (c) Concept Unique Identifier (CUI) reuse. While there is a term label similarity of approximately 14.4% of the total terms, we observed that most ontologies reuse considerably fewer than 5% of their terms from a concise set of a few core ontologies. We developed an interactive visualization to explore reuse dependencies among biomedical ontologies. Moreover, we identified a set of patterns that indicate ontology developers did intend to reuse terms from other ontologies, but they were using different and sometimes incorrect representations. Our results suggest the value of semi-automated tools that augment term reuse in the ontology engineering process through personalized recommendations.

1 INTRODUCTION

Ontologies have been used in biomedical research for different purposes: knowledge management, semantic search, data annotation, data integration, exchange, decision support and reasoning (Bodenreider, 2008; Rubin *et al.*, 2008). Biomedical ontologies range drastically in their size, coverage of a domain, and in their level of adoption. It is only natural that biomedical ontologies will overlap to a certain degree, as they sometimes need to represent common parts of a domain, or different domains that have shared terms.

Several large biomedical efforts deal in different ways with managing the overlap of ontologies and reuse. For example, the Open Biological and Biomedical Ontologies (OBO) Foundry (Smith *et al.*, 2007) aims to create a set of “orthogonal” ontologies, such that each term is defined only in one ontology, and is referred using its Internationalised Resource Identifier (IRI) in other ontologies. The OBO ontologies use the *xref* mechanism to create references between terms in different ontologies (OBOFoundry, 2011). To support the interoperability across different biomedical ontologies and terminologies, the Unified Medical Language System–UMLS (Bodenreider, 2004) uses the notion of a Concept Unique Identifier (CUI) to map terms with similar meaning in different terminologies.

All ontology development methodologies strongly encourage reuse while building new ontologies, be it at the level of an ontology, or at the level of the terms (Corcho *et al.*, 2003; Alexander, 2006). Reuse has some directly apparent advantages, such as, developing a unified theory of biomedicine, semantic interoperability and reducing engineering costs (since reuse avoids rebuilding existing ontology structures). For example, the 11th revision of the International Classification of Diseases (ICD-11) reuses terms from other established ontologies, such as SNOMED

CT, to spare the effort in creating already existing quality content (e.g., the anatomy taxonomy), to increase their interoperability, and to support its use in electronic health records (Tudorache *et al.*, 2010). Another benefit of the ontology term reuse is that it enables federated search engines to query multiple, heterogeneous knowledge sources, structured using these ontologies, and eliminates the need for extensive ontology alignment efforts (Kamdar *et al.*, 2014).

For the purpose of this work, we define as **term reuse** the situation in which the same term is present in two or more ontologies either by direct use of the same identifier, or via explicit references and mappings. We define as **term overlap** the situation in which the term labels in two or more ontologies are lexically similar (see Section 3). We further classify the reuse: (1) *Reuse of an ontology*, through the means of the import mechanism available in OWL (W3C, 2012), meaning that the entire source ontology is imported into the target ontology; and (2) *Reuse of terms* from one source ontology into another. In many cases, experts reuse not only one term from one ontology, but rather subsets of terms from multiple ontologies (e.g., subtrees).

The goal of this work is to investigate the level of reuse and overlap among biomedical ontologies. We harvested the ontologies from BioPortal (Whetzel *et al.*, 2011), an open content repository of biomedical ontologies and terminologies, and ran several analyses that show not only the level of reuse, but also how the reuse occurs.

The contributions of this work are threefold:

1. A set of descriptive statistics for the level of reuse in biomedical ontologies,
2. An interactive visualization technique for displaying the reuse dependencies among biomedical ontologies,
3. A discussion on the state and challenges of reuse in biomedical ontologies.

2 RELATED WORK

Through a set of use cases in bio-medicine and eRecruitment, Bontas *et al.* (2005) emphasise the need for more pragmatic methods and semi-automated tools that allow developers to exploit the vocabulary of domain-specific source ontologies for reuse. Matentzoglou *et al.* (2013) provide a method to analyze the overlap between ontologies in automatically-generated random snapshots of the OWL Web. Ontologies with 90% overlap or containment relations were considered similar. Ghazvinian *et al.* (2011) describe an approach to determine the level of orthogonality and *term overlap* (term label similarity) in OBO Foundry member and candidate ontologies. Their analysis over a period of two years indicated that, while the OBO Foundry has made significant progress towards achieving “orthogonality”, *term overlap* between

ontologies has remained consistent. Poveda Villalón *et al.* (2012) analyze the landscape of reuse in the ontologies used in Linked Open Data (LOD), and find that over 40% of the terms are reused from other vocabularies, 67% of which are reused by imports, and the rest by referencing the term IRI. Ontology modularisation techniques (i.e., extracting parts of an ontology using some structural or logical properties) are also an important factor in supporting reuse. Comprehensive studies of existing modularization techniques have also been undertaken (d'Aquin *et al.*, 2009; Pathak *et al.*, 2009).

There are only a few tools that support term reuse in biomedical ontologies. OntoFox (Xiang *et al.*, 2010) is a Web-based application that allows users to retrieve terms, selected properties, and annotations from the source ontologies, using MIREOT principles (Courtot *et al.*, 2011). The BioPortal Import Plugin (Nair, 2014), MIREOT Protégé Plugin (Hanna *et al.*, 2012) and DOG4DAG (Wächter *et al.*, 2011) are provided as extensions to the Protégé ontology editor (Noy *et al.*, 2001) to allow the import of terms, their properties, and even class subtrees from BioPortal.

3 METHODS

We obtained a triplestore dump of the BioPortal ontologies in N-triples format as of January 1, 2015, which contained 377 distinct biomedical ontologies. This dump does not contain some ontologies that were deprecated or merged with existing ontologies, or those that were added to BioPortal after January 1, 2015. These ontologies include eight OBO Foundry member ontologies (GO, CHEBI, PATO, OBI, ZFA, XAO, PR and PO) and 31 UMLS Terminologies (SNOMED CT, ICD-9, etc.). To conduct our analysis, we identified three constructs that cover reuse in BioPortal ontologies:

1. **Explicit reuse construct:** The IRIs of the terms in different ontologies are exactly the same.
2. **xref construct:** One term contains a reference to the other term IRI using the *xref* predicate.
3. **UMLS CUI construct:** Two BioPortal-defined term IRIs are mapped to the same Concept Unique Identifier.¹

By iterating over all the asserted axioms in each of these 377 ontologies, we extracted all the class term IRIs, their labels, *xref* links and UMLS CUI mappings, when available. From the 5,718,276 class terms, we used the above three constructs to extract the set of terms that satisfy any of the three reuse criteria. The *xref* axioms were further filtered to separate those that assert equivalence between the connected entities (e.g., `CL:0000066`, `CARD:0000077` and `FMA:66768` all refer to 'epithelial cell'), from those that were either references to resources in external databases like PubMed, or entities that were semantically treated as *genus-differentia* definitions, as defined in the OBO Foundry (2011).

For the first two reuse types (*Explicit* and *xref*),² we identified the source ontology for each term by converting each term IRI to lowercase and using *RegExp*

filters constructed from ontology namespaces and common identifier patterns. A heuristic approach was used to determine the source ontology, by first checking only the current ontologies that share this term. We found some instances where the source is not determined in the first step. For example, `NCIT:Cerebral_Vein` is reused by Sage Bionetworks Synapse Ontology (SYN) and Cigarette Smoke Exposure Ontology (CSEO). However, this term is replaced by `NCIT:C53037` in the current version of NCI Thesaurus, and the original term is not present. Hence, as a second step, we extended our search to include all the ontologies. This two-step approach also deals with the conditions when an ontology acronym ('PR') is present in a term IRI (e.g., `Protein`) but is not necessarily in the source ontology.

We normalised the term labels by converting them to lowercase and removing all non-alphanumeric characters. We performed naïve string matching on the term labels to determine the potential *term overlap*.

We calculated three statistics:

1. The percentage of terms explicitly reused or *xref*-linked by an ontology, and the total number of ontologies these terms are reused from (on *Explicit* and *xref* constructs),
2. The percentage of terms and the total number of ontologies that are reused explicitly, or *xref*-linked, from an ontology (on *Explicit* and *xref* constructs).
3. The reuse between all distinct pairs of ontologies (on *Explicit*, *xref*, and *UMLS CUI* constructs).

Using these statistics, we determined which ontologies reused the maximum number terms from other ontologies, and also those ontologies whose terms were reused the most. We calculated the gap between term overlap and term reuse by subtracting the matched labels of reused terms.

To determine the level of import at the level of an ontology, we used the explicit occurrence of the `owl:imports` in the ontology files. However, this method did not account for the cases in which the imported ontologies were already merged into the importing ontology, as is the case in BioPortal. Therefore, we established an empirical threshold of 35% on the number of terms that were reused with respect to the total number of classes in the source ontology, above which we would consider the term reuse as a *reuse of the entire ontology*. As determined from reuse statistics, this threshold would allow us to consider term reuse from older versions of source ontology as ontology reuse.

During the development of an ontology for a specific domain, it is beneficial for the ontology engineers to have an idea regarding the set of ontologies whose terms were reused from by other related ontologies. Hence, we developed an interactive force-directed network visualization to represent the ontology pairs derived from the third statistic to explore the reuse dependencies among biomedical ontologies.

4 RESULTS

Explicit Reuse First, we investigated the **reuse at the level of an ontology** by the means of `owl:imports` mechanism and the 35% threshold method. The top 10 of the most imported ontologies are shown in Table 1.

¹ This was only checked for UMLS terminologies.

² UMLS CUI reuse was excluded, as we could not identify the source ontology for a CUI.

Ontology Name	#
(BFO) Basic Formal Ontology	59
(STY) Semantic Types Ontology	29
(PATO) Phenotypic Quality Ontology	10
(IAO) Information Artifact Ontology	9
(UO) Units of Measurement Ontology	5
(CARO) Common Anatomy Reference Ontology	4
(ONL-MSA) OntoNeuroLOG - Mental State Assessment	3
(ORDO) Orphanet Rare Disease Ontology	2
(GO) Gene Ontology	2
(BP) BioPAX Ontology of Biological Pathways	3

Table 1. Most imported ontologies (Reuse of an ontology). (#) indicates number of ontologies importing the specified ontology.

Second, we investigated the explicit **reuse of individual terms**. Of the 5,718,276 class terms that we extracted from the 377 BioPortal ontologies, we found 175,347 terms (3.1%) were explicitly shared among more than two ontologies using the same IRIs. We found the source ontology for all but 37 terms, which were primarily upper-level, abstract terms, whose ontologies were not present in BioPortal (e.g., `owl:Thing` and `time#datetimedescription`). After removing the terms that come from imported ontologies that were merged (term reuse > 35% threshold), we were left with only 59,618 terms (1.1%) actually reused.

xref Reuse We found a total of 4,370,350 *xref* axioms across all the BioPortal ontologies. After extracting *xrefs*, which assert equivalence between BioPortal ontology terms, we found 171,069 ‘outlinking’ terms (3.9%) *xref*-linked to 386,442 ‘inlinking’ terms (8.84%).

We also tried to understand how the explicit- and *xref* reuse is spread across different ontologies. **Figure 1** shows histograms of the percentage of terms reused by different ontologies. We can see that most ontologies reuse or *xref*-link less than 5% of their total terms. There were at least 150 ontologies which did not reuse a single term from other ontologies. We also observe that there are 20 ontologies that exhibit a reuse between 95% to 100% of their total terms. These ontologies are developed by reusing combinations of multiple ontologies (e.g., CCONT reuses terms from EFO, NCBITAXON, ORDO, and 19 other ontologies).

The top 10 ontologies that reuse their terms from the maximum number of ontologies, and those whose terms are reused the most, are shown in Tables 2 and 3, respectively. The columns indicate the percentage (%) of total terms explicitly reused or *xref*-linked from/by the number of ontologies (#). For example, current version of NIFSTD explicitly reuses 89.6% of its total terms from 42 different ontologies, and 95.2% and 3.7% of the total terms in the current version of GO are reused and *xref*-linked by 74 and 37 ontologies respectively. The top 10 terms, which are not upper ontology terms (e.g., from BFO or IAO) and are explicitly reused the most, are shown in Table 4.

UMLS CUI Reuse Using our third construct, we found 236,460 Concept Unique Identifiers (CUIs), which are mapped to more than two terms in UMLS terminologies. Some of the most mapped CUI terms are: **Neoplasms** (C0027651) and **Diabetes mellitus** (C0011849) appearing in 18 ontologies, **Schizophrenia** (C0036341) and

Ontology (<i>explicit</i>)	% Reused	#	Ontology (<i>xref</i>)	% <i>xref</i> -linked	#
NIFSTD	89.6	42	UBERON	72.2	37
HUPSON	55.8	32	CL	14.2	21
OBI_BCGO	97.9	25	TMO	17.3	21
IDOMAL	43.5	24	HPIO	53.7	16
IDODEN	29.1	23	DOID	90.8	13
OBI	19.1	22	TRAK	23.8	10
CCONT	98.8	22	GO	0.76	9
EFO	70.1	21	HP	11.8	8
CLO	7.2	19	DERMO	25.7	7
IDOBRU	43.27	19	EFO	0.76	6

Table 2. Ontologies that reuse the maximum number of terms from other ontologies - Percentage (%) of the total number of terms reused from the total number of ontologies (#).

Ontology (<i>explicit</i>)	% Reusing	#	Ontology (<i>xref</i>)	% <i>xref</i> -linking	#
BFO	259	81	GO	3.7	24
GO	95.2	74	CHEBI	3.2	16
IAO	72.8	55	CARO	572	16
OBI	43.1	51	MESH	2.3	11
PATO	190	45	PATO	23.6	10
CHEBI	54.2	37	FMA	14.0	10
CL	15.4	36	NCIT	6.6	10
NCBITAXON	0.30	30	CL	18.9	9
STY	100	29	NCBITAXON	19.9	8
UO	136	27	SO	5.0	8

Table 3. Ontologies whose terms are reused most by other ontologies - Percentage (%) of the total number of terms in the current version reused by the total number of ontologies (#).

Term IRI	Term Label	#Reusing Ontologies
GO:0008150	biological_process	33
OBI:0000011	planned process	31
OBI:0100026	organism	28
CHEBI:23367	molecular entity	26
NCBITaxon:9606	Homo sapiens	26
PATO:0001241	physical object quality	24
GO:0005575	cellular_component	23
PATO:0001995	organismal quality	23
PATO:0000001	Quality	22
NCBITaxon:10239	Virus	20

Table 4. Top 10 terms that are reused by maximum ontologies

Leukemia (C0023418) appearing in 17 ontologies. The full list is available online (see link at the end of section). **Figure 2** shows the percentage of CUI terms shared by each UMLS terminology with other terminologies. It is noteworthy to see some of the popular UMLS terminologies such as ICD10CM, LOINC, HL7 and MESH to be composed primarily of unshared, unique terms.

Overlap Executing normalised string matching on the term labels, we found a *term overlap* of 823,621 shared term labels (14.4%). On removing those terms that were already explicitly reused using the same term IRI, we reduced the list to 752,176 labels (13.2%). On removing those terms which were mapped to the same UMLS CUI, we further reduced the list to 617,509 labels (10.8%). On extracting

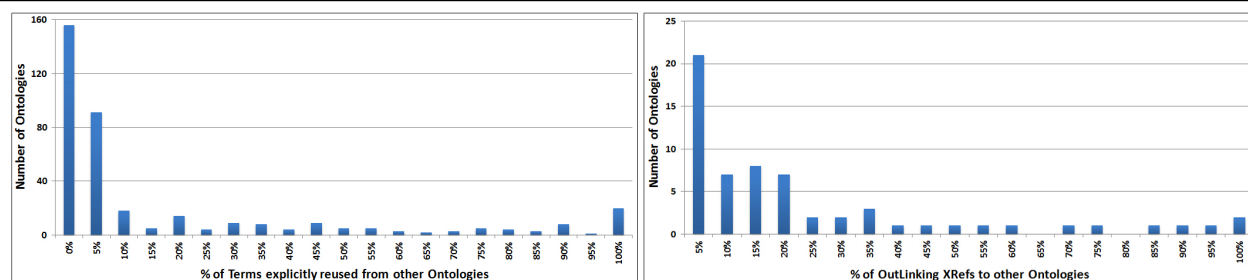


Fig. 1. Histograms depicting the first statistic: a) Percentage (%) of terms explicitly reused or b) *xref*-linked by an ontology

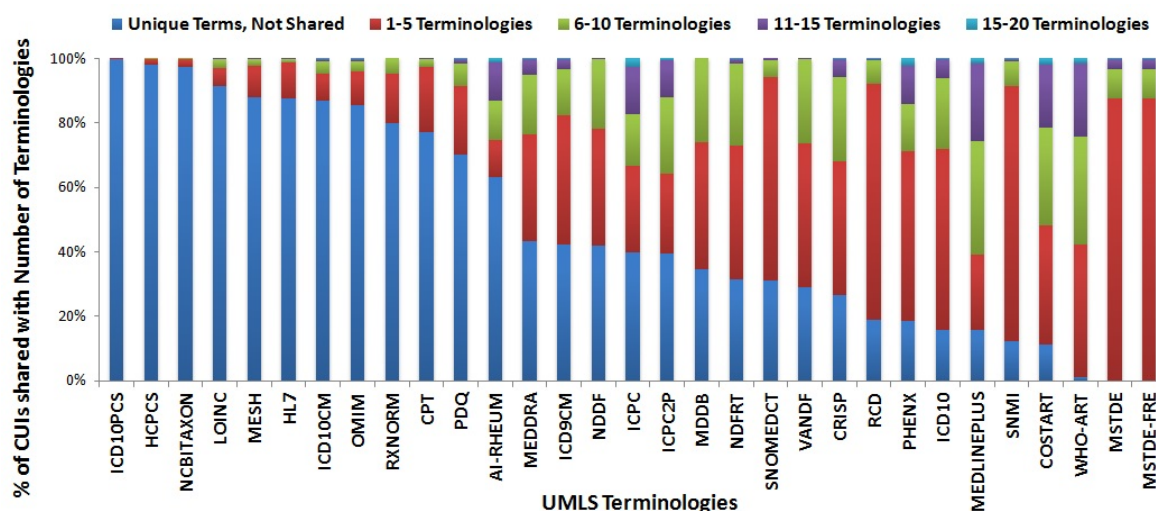


Fig. 2. Percentage (%) of CUIs in each UMLS terminology shared with other terminologies

the resource identifier from each term IRI, we removed those terms which had almost similar term IRIs (having the same identifier and source ontology, but a different or incorrect representation), and the list reduced to 93,650 term labels (1.6%). The last step does not represent actual reuse between ontologies, but rather that ontology developers showed an intention to reuse terms, but used different and sometimes incorrect term representations (discussed below).

Force-directed network visualization We developed an interactive force-directed network visualization, using the third statistic, where the ontologies form the nodes of the networks and the edges connecting them indicate the extent of term reuse between them. The nodes are colored according to the group under which the ontology falls, and the size of the nodes depends on the total number of terms in the current version of the ontology. The thickness of the edge is proportional to the total number of terms shared between the connected ontologies, and the colour varies according to the construct. The graph can be constrained by hovering over any node, to display only the directly related nodes. The interactive version can be accessed at: <http://stanford.edu/~maulikrk/apps/OntologyReuse/>.

The detailed results are available at <http://stanford.edu/~maulikrk/data/OntologyReuse/>.

5 DISCUSSION

Ghazvinian *et al.* (2011) outlined the consistent term overlap, yet minimum term reuse, in OBO Foundry ontologies, and commented on the limitations and challenges to achieve

“orthogonality”. Five years later, evaluating term reuse over the entire continuum of biomedical ontologies (including UMLS terminologies), we see that we are still very far from achieving desirable term reuse. Most ontologies exhibit considerably less than 5% reuse or no reuse through any constructs, and generally reuse terms from only a small set of ontologies. Table 2 lists many of the OBO Foundry member ontologies. The OBO Foundry mandates reuse by candidate ontologies from the member ontologies under its orthogonality aim. However, there is still substantial *term overlap* present among biomedical ontologies, including OBO Foundry ontologies. *Term overlap*, in itself, is not a good indicator of potential term reuse, as there may be terms in different ontologies which are lexically similar, but represent different concepts (e.g., similar anatomical concepts between Zebrafish Anatomy (ZFA) and Xenopus Anatomy (XAO)), and lexically-different terms may represent the same concept (e.g., *myocardium* and *cardiac muscle*). Hence rigorous methods to detect contextual term overlap are required.

By examining the terms that shared the same labels, we found various IRI patterns that could indicate that the ontology developers showed the intention to reuse terms (same identifiers and source ontologies). These patterns were not considered as term reuse as the IRIs used different representations for the same terms, and no explicit CUI or *xref* mappings were found. Hence, the advantages of term reuse can not be experienced. On using the right IRI representation, the term overlap could reduce substantially. We describe the three most prevalent patterns below.

Different versions: SAO and SOPHARM reuse terms from BFO version 1.0, whereas the majority of other ontologies reuse the corresponding terms found in version 1.1. As mentioned in Section 3, CSEO and SYN reuse terms from an older version of NCI Thesaurus. For example, we found NCIT:Cerebral_Vein instead of the recent NCIT:C53037.

Different notations: Terms reused from FMA were referenced in multiples ontologies using different notations without consistency or interlinks. For example, OBO:FMA_31396 is reused as OBO:owlapi/fma#FMA_31396, OBO:owl/FMA#FMA_31396, and even with the entire label OBO:fma#Cartilage_of_inferior_surface_of_posterolateral_part.

Different namespaces: Different ontologies tend to use completely different namespaces for the source ontology. For example, RH-MESH uses <http://phenomebrowser.net/ontologies/mesh/mesh.owl>, while most other ontologies reuse <http://purl.bioontology.org/ontology/MESH>. We also found reuse of SNOMED CT terms with two distinct namespaces: <http://ihtsdo.org/snomedct/clinicalFinding> and <http://purl.bioontology.org/ontology/SNOMEDCT>.

There are direct (semantic interoperability, cost reduction) and indirect (EHR mining, query federation) advantages of term reuse. In the Linked Open Dataspace, newer, collaborative efforts, such as Bio2RDF (Callahan *et al.*, 2013), provide strict guidelines for the representation of concept identifiers while publishing data as RDF. ProtégéLov (Garcia-Santa *et al.*, 2015) allows reuse of terms directly from the Linked Open Vocabularies repository using `owl:equivalentClass` and `rdf:subClassOf` axioms.

Our analysis indicate that while ontology developers may exhibit an intention for term reuse, the lack of guidelines and semi-automated tools for ontology term reuse seem to hinder these goals. Our visualization of reuse dependencies could guide developers to reuse terms in their own ontology based on the structure of ontologies in related domains. Identifying reuse patterns and providing personalized recommendations during the development phase could help increase term reuse and reduce term overlap. Incorporating a reuse module in ontology editing tools could also keep developers updated when the representation of the source term changes.

6 CONCLUSION

We analyzed the extent of term reuse and overlap in 377 biomedical ontologies from BioPortal along three reuse constructs: explicit reuse, *xref* reuse, and CUI reuse. Despite the considerable level of overlap (14.4%), there is very little reuse (< 5%) among biomedical ontologies, both at the level of an ontology and at the level of individual terms. We developed a force-based visualization that helps users to understand the reuse dependencies across different ontologies. We also identified error patterns in applying reuse that we discovered in our empirical analysis. Our future work includes research on identifying reuse patterns in an empirical way, and building a recommendation module for the Protégé toolset that would suggest terms that have been reused together with existing terms. Our strong belief is that better guidelines and tool support will enhance the reuse among biomedical ontologies.

ACKNOWLEDGMENTS

The authors acknowledge Manuel Salvadores Olaizola for providing a triplestore dump of BioPortal ontologies. This work is supported in part by grants GM086587 and GM103316 from the US National Institutes of Health.

REFERENCES

- Alexander, C. Y. (2006). Methods in biomedical ontology. *Journal of biomedical informatics*, **39**(3), 252–266.
- Bodenreider, O. (2004). The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic acids research*, **32**(suppl 1), D267–D270.
- Bodenreider, O. (2008). Biomedical ontologies in action: role in knowledge management, data integration and decision support. *Yearbook of medical informatics*, page 67.
- Bontas, E. P. *et al.* (2005). Case studies on ontology reuse. In *Proceedings of the IKNOW05*, volume 74.
- Callahan, A. *et al.* (2013). Bio2RDF release 2: Improved coverage, interoperability and provenance of life science linked data. In *The Semantic Web: Semantics and Big Data*, pages 200–212. Springer.
- Corcho, O. *et al.* (2003). Methodologies, tools and languages for building ontologies. Where is their meeting point? *Data & knowledge engineering*, **46**(1), 41–64.
- Courtot, M. *et al.* (2011). MIREOT: The minimum information to reference an external ontology term. *Applied Ontology*, **6**(1), 23–33.
- d’Aquin, M. *et al.* (2009). Criteria and evaluation for ontology modularization techniques. In *Modular ontologies*, pages 67–89. Springer.
- Garcia-Santa, N. *et al.* (2015). *Protege LOV Plugin*. <http://goo.gl/9fmTf7> (accessed March 05, 2015).
- Ghazvinian, A. *et al.* (2011). How orthogonal are the OBO foundry ontologies? *J. Biomedical Semantics*, **2**(S-2), S2.
- Hanna, J. *et al.* (2012). Simplifying MIREOT: a MIREOT protégé plugin. In *The Semantic Web – ISWC 2012*.
- Kamdar, M. R. *et al.* (2014). ReVeLD: A user-driven domain-specific interactive search platform for biomedical research. *Journal of biomedical informatics*, **47**, 112–130.
- Matentzoglou, N. *et al.* (2013). A snapshot of the OWL web. In *The Semantic Web – ISWC 2013*, pages 331–346. Springer.
- Nair, J. (2014). *BioPortal Import Plugin*. <http://goo.gl/LL75TR> (accessed March 01, 2015).
- Noy, N. F. *et al.* (2001). Creating semantic web contents with protege-2000. *IEEE intelligent systems*, **16**(2), 60–71.
- OBOFoundry (2011). *Inter-ontology Links*. <http://goo.gl/OSrSjP> (accessed March 01, 2015).
- Pathak, J. *et al.* (2009). Survey of modular ontology techniques and their applications in the biomedical domain. *Integrated computer-aided engineering*, **16**(3), 225–242.
- Poveda Villalón, M. *et al.* (2012). The landscape of ontology reuse in linked data. In *Proceedings of OEDW 2012*. Informatica.
- Rubin, D. L. *et al.* (2008). Biomedical ontologies: a functional perspective. *Briefings in bioinformatics*, **9**(1), 75–90.
- Smith, B. *et al.* (2007). The OBO Foundry: coordinated evolution of ontologies to support biomedical data integration. *Nature biotechnology*, **25**(11), 1251–1255.
- Tudorache, T. *et al.* (2010). Ontology development for the masses: creating ICD-11 in WebProtégé. In *Knowledge Engineering and Management by the Masses*, pages 74–89. Springer.
- W3C (2012). *OWL 2 Web Ontology Language Document Overview*. <http://www.w3.org/TR/owl2-overview/> (accessed March 01, 2015).
- Wächter, T. *et al.* (2011). DOG4DAG: semi-automated ontology generation in obo-edit and protégé. In *Proceedings of SWAT4LS 2011*, pages 119–120. ACM.
- Whetzel, P. L. *et al.* (2011). BioPortal: enhanced functionality via new Web services from the National Center for Biomedical Ontology to access and use ontologies in software applications. *Nucleic acids research*, **39**(suppl 2), W541–W545.
- Xiang, Z. *et al.* (2010). OntoFox: web-based support for ontology reuse. *BMC research notes*, **3**(1), 175.