

CONTEXT-SPECIFIC ONTOLOGY INTEGRATION*

KSHITIJ MARWAH

*Department of Computer Science and Engineering, Indian Institute of Technology
New Delhi, Delhi 110016, India.
Email: ksm@mit.edu*

DUSTIN KATZIN

*Department of Physics, Massachusetts Institute of Technology
Department of Mathematics, Massachusetts Institute of Technology
Cambridge, Massachusetts 02139, USA.
Email: drkatzin@mit.edu*

GIL ALTEROVITZ†

*Children's Hospital Informatics Program, Harvard-MIT Division of Health Sciences and Technology
Harvard Medical School, Boston, Massachusetts 02115, USA.
Partners Healthcare Center for Personalized Genetic Medicine.
Boston, Massachusetts 02115, USA.
Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology
Cambridge, Massachusetts 02139, USA.
Email: gil_alterovitz@hms.harvard.edu*

Today, ontologies have become critical tools of biomedical research, providing an efficient framework for structuring and organizing scientific information. Integration of many different biomedical ontologies into a comprehensive landscape of biomedical knowledge would enable researchers to generate new hypotheses and identify new avenues of investigation. Here we introduce a principled computational framework for automated discovery of functional links amongst ontologies. We have developed a novel technique, Lexical Mapping, for deducing context specific functional links, by leveraging over disparate free-text literature resources. We start by searching for ontological terms over the corpus and caching the terms and annotations using a b-tree based reverse index. Next, we use a pre-computed transitive closure of the ontology graph to propagate the annotations up the hierarchy, and thus augment it to the index. To infer context-specific links, we score the model of dependency linking two terms under the given context against their model of independence, further penalizing it using the bayesian information criterion to get the bayes factor. We finally identify context-specific linked terms as those having a bayes factor greater than twenty ($p < 0.01$). To scale our algorithm over large ontologies, we develop a heuristic pruning technique, using a depth first branch and bound algorithm to exponentially reduce running time with marginal loss in the quantity of inferred links. To further augment our inferred links, we compose the existing links to relate different ontologies already integrated into a given ontology under the same context. Using the functional meanings of biomedical terms to make connections, this approach is fundamentally different from its predecessors, which primarily relied on connecting the lexical structures of term names. We have applied this method to translationalize Gene Ontology to all other ontologies available at National Center of Biomedical Ontologies (NCBO) Biportal, under the context of the Human Disease ontology. This is the first time that all ontologies

* Work partially supported by grants 1K99LM009826 and 5T15LM007092 of the National Library of Medicine (NLM/NIH), as well as grants 2P41HG02273, 1R01HG003354, and 1R01HG004836 of the National Human Genome Research Institute (NHGRI/NIH)

† To whom correspondence should be addressed

‡ Deceased

have been included in one map. We have validated our inferred links manually by supporting published literature and by a domain expert. We believe, in addition to broadening the scope of hypotheses for researchers, our work can also be used to explore the continuum of relationships amongst ontologies to guide various biological experiments.

1. Introduction

Every year, over 400,000 new articles reportedly enter biomedical literature [1]. This staggering growth of biomedical findings has created an unprecedented corpus of knowledge that is impossible to explore with traditional means of literature consultation and database searches. This information overload has motivated the development of structured information repositories that organize biomedical findings according to hierarchical ontologies.

Ontologies are currently at the heart of two major complementary activities in biomedical research. Firstly, communities of researchers create and maintain these ontologies to represent different types of entities and relations in different domains of biomedicine. Secondly, biomedical experimentalists use ontologies to annotate data enabling their information to be integrated with other researchers data, and permitting cross-species analyses through experimental data annotations. This activity is greatly intensified by the development of high-throughput experimental platforms such as gene expression microarrays [2], SNP microarrays [3], and next generation sequencing platforms [4].

The rise of such ontological organization has created a new problem, the proliferation of disparate and seemingly unrelated biomedical ontologies available to researchers. The National Center of Biomedical Ontology (NCBO) Bioportal [5], provides over 200 such ontologies to researchers, which is of great assistance to them. Often, researchers will need to use multiple ontologies to annotate their data, but which ontologies to use, and how they relate to one another is generally unclear. What is needed is the integration of the various ontologies in a principled fashion, a “grand unification” of biological terms. It has been established [6], that the integration of these available biomedical ontologies will have a tremendous impact on the advance of the biomedical sciences. These integrated ontologies would provide a complete basis of biomedical knowledge representation, and would act as a foundation for inference on new biomedical data. Furthermore, a quantitative approach to integration would make the navigation of the complex space of ontologies more amenable to biomedical researchers, by offering them guidance to numerous links amongst discrete ontologies, thus making the discovery process faster and more efficient.

To date, approaches to mapping and integrating ontologies have relied on discovering links between semantically similar terms across ontologies [7]. Such an approach can relate terms with similar meanings, but would not deduce, say, relationships between seemingly disparate functional spaces such as diseases, drugs and anatomy. Other approaches to infer such types of links, use standard means of manual curation, which is a tedious and labor intensive task, unable to scale up to the current size of biomedical ontologies and keep up with their growth rate.

Here we propose a novel computational and methodological framework for context specific integration of biomedical ontologies using free-text literature analysis. We model context

specificity using another ontology, and derive context-dependent functional links between ontological concepts occurring as terms in free-text literature. We conduct our search for ontological concepts in free-text biomedical literature, and efficiently caching the required statistics using a b-tree based reverse index. We then leverage upon these cached statistics to compute the penalized likelihood of the model of dependency and the model of independence by applying the well-known bayesian information criterion (BIC) [8], over a context-sensitive model scoring function. Using the BIC, we then deduce the most likely model explaining the observed data, thus classifying the context-specific relation between terms as related or not related. We also rank our links according to a principled bayes factor metric with direct correspondence to the strength of our findings.

Due to large scale of ontologies involved, and the complexity of the search space, a “brute force” computational approach would not scale using these techniques. To help circumvent this problem, we propose a depth first branch and bound heuristic pruning technique to help us prune away subgraphs of the ontologies which would not yield significant functional links. We show how such techniques result in exponential reduction in running time with marginal loss in quantity of links. To further augment our inferred links, we apply mapping composition to yield mappings amongst ontologies integrated via another ontology. This allows us to use the existing mappings to generate links amongst ontologies that have not been directly mapped.

We believe that such a methodological approach would turn available machine-processable ontologies into a single landscape of integrated biomedical concepts and annotations. This would thus enable researchers to bear on each single finding, the entire power of established biomedical knowledge.

2. Methods

In all, 200 ontologies from the National Center of Biomedical Ontology’s Bioportal interface were obtained. To enable us compute the likelihood of dependency amongst ontology concepts we develop a pipeline similar to the NCBO Annotator [10], but targeted towards efficient frequent counting and information retrieval.

2.1. *Caching Sufficient Statistics*

We gather raw free-text literature from disparate sources, and drive our concept search by finding exact matches of ontology terms. We use the MGREP [11], concept recognition tool to efficiently find occurrence of concepts in published literature and thus annotate the documents with those concepts. This allows us to leverage on a consolidated vocabulary (of about 4 million ontology concepts) to temper the problem of missing synonyms and term permutations.

We also use a pre-computed index containing the transitive closure of ontology terms for semantically expanding the annotations, propagating them up the hierarchy of the ontology. The document annotations and the concepts are reverse indexed using a disk based b-tree data structure, an approach commonly used in information retrieval domains.

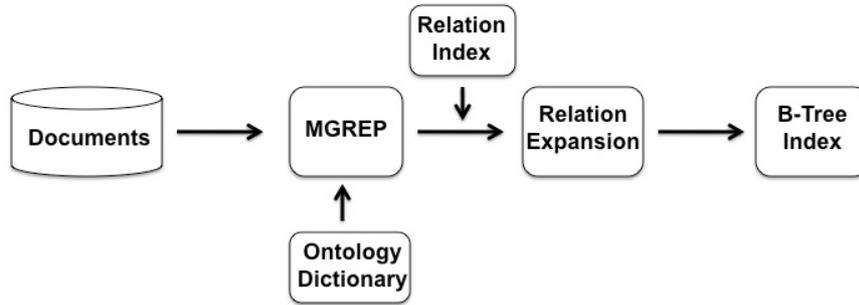


Figure 1. Pipeline used for caching sufficient statistics for model scoring.

We use Lucene [12], an open source high-powered information retrieval engine to create and store the b-tree structure. To answer conjunctive queries for efficient counting, we use a bitmap hash based filter over the stored index. Our analysis over disparate corpora shows linear scalability of the data structure in terms of time and space on account of new corpora and new ontologies. The query time for counting was found to scale logarithmically with the addition of new ontologies. Further increase in efficiency was observed, by cleverly caching filters when computing frequent counts.

2.2. Lexical Mapping Algorithm

For computing context dependent links between ontology terms, we have developed a novel technique called Lexical Mapping, which relies on the statistical analysis of literature. Lexical mapping uses the observed co-occurrence of terms in the literature to infer the relationship between two terms A and B in the context of the ontology term C.

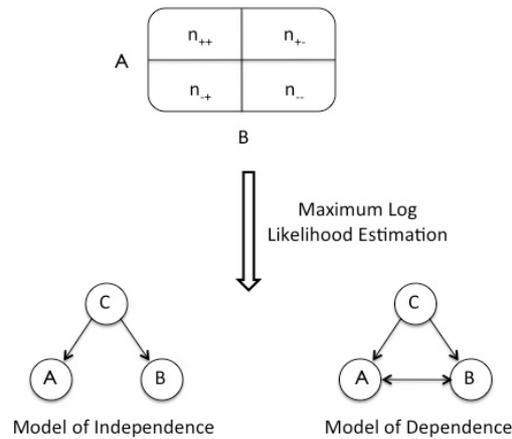


Figure 2. 2 x 2 contingency table to test relationship between two ontology terms A and B under C.

To do so it builds a contingency table like the one in Figure 2, collecting the frequencies of co-occurrence of the two terms in the literature, a 2 x 2 table where n_{++} is the number of cases in which the two terms appear together, n_{+-} is the number of cases in which A appears but B doesn't, n_{-+} is the number of cases in which B appears but A doesn't, n_{--} is the number of cases in which neither appear in the context of term C.

Our method uses the bayesian information criterion to compute the penalized likelihood of dependency $A \Leftrightarrow B \mid C$ (where two terms are related) and the model of independence $A \Uparrow B \mid C$ (where the two terms are unrelated) as

$$BIC = -2MLL + k \log(N) \quad (1)$$

where N is the number of observations, k is the number of parameters of the model, and MLL is the marginal log likelihood of the model.

The marginal log likelihood for the model of dependency is

$$\begin{aligned} MLL(A \Leftrightarrow B \mid C) = & [\ln(\Gamma(\alpha)) - \ln(\Gamma(\alpha + n))] \\ & + [\ln(\Gamma(\alpha_k + n_{++})) - \ln(\Gamma(\alpha_k))] \\ & + [\ln(\Gamma(\alpha_k + n_{+-})) - \ln(\Gamma(\alpha_k))] \\ & + [\ln(\Gamma(\alpha_k + n_{-+})) - \ln(\Gamma(\alpha_k))] \\ & + [\ln(\Gamma(\alpha_k + n_{--})) - \ln(\Gamma(\alpha_k))] \end{aligned} \quad (2)$$

while the marginal log likelihood for the model of independence is

$$\begin{aligned} MLL(A \Uparrow B \mid C) = & [\ln(\Gamma(\alpha)) - \ln(\Gamma(\alpha + n))] \\ & + [\ln(\Gamma(\alpha_k + n_{++} + n_{+-})) - \ln(\Gamma(\alpha_k))] \\ & + [\ln(\Gamma(\alpha_k + n_{-+} + n_{--})) - \ln(\Gamma(\alpha_k))] \\ & + [\ln(\Gamma(\alpha)) - \ln(\Gamma(\alpha + n))] \\ & + [\ln(\Gamma(\alpha_k + n_{-+} + n_{++})) - \ln(\Gamma(\alpha_k))] \\ & + [\ln(\Gamma(\alpha_k + n_{--} + n_{+-})) - \ln(\Gamma(\alpha_k))] \end{aligned} \quad (3)$$

where Γ is the gamma function, n_{++} , n_{+-} , n_{-+} , n_{--} are the co-occurrence frequencies as described above, α is the prior precision and, α_k is the prior precision per term that is $\alpha/|T|$, where $|T|$ is the number of terms in the dependency, in our particular case 2. In our case we use $\alpha = 4$ for 2×2 tables, so that for the initial prior precision we put 1 in each cell, maintaining the uniformity of the distribution and the lowest possible precision, so as to minimize bias on the precision.

By plugging the marginal log likelihood into equation (1), we obtain respectively the penalized likelihood for the model of dependency $BIC(A \Leftrightarrow B \mid C)$, where the two terms are linked, and the model of independence $BIC(A \Uparrow B \mid C)$, where the two terms are not linked. The final score is the bayes factor

$$Score = BIC(A \Leftrightarrow B \mid C) / BIC(A \Uparrow B \mid C) \quad (4)$$

which estimates how many times the model linking term A and B in the context of C is more likely than the model in which the terms are not related.

We use our cached data structure to efficiently count the co-occurrence frequencies, for computing the bayes factor. Context-dependent functional links are then selected as the ones having bayes factor greater than 20 ($p < 0.01$).

2.3. Heuristic Pruning Using A Depth First Branch And Bound Algorithm

To apply the Lexical Mapping algorithm, we in the worst case would have to compute for all possible triples of terms representing the ontologies. Such an approach though would work for small ontologies, does not scale up to large ontologies even with the efficient caching data structure. We apply a depth first branch and bound algorithm and prune away ontology subgraphs where the likelihood of finding functional links is extremely low. We use the bayes factor as a scoring cue to find such subgraphs.

Theorem 1. If the bayes factor for a ontology concept A, mapped to another ontology concept B under the context of C is less than a given threshold ϵ , then the likelihood of finding a map amongst a major fraction of A's children, with the concept B under C also decreases.

Proof. We use the fact that in an ontology, any instance of a specific concept is also an instance of a more general super concept. This implies that while propagating instances up from A's children to the concept A, if enough evidence of a link was found between B and A's children under C, that evidence would have propagated up the hierarchy, thus linking A to B under C. This further implies that, for A to be mapped to B under C, a major fraction of A's children should also be linked to B under C, thus propagating enough evidence for A to be linked. By transposition, if the map of A to B under C is very less likely, the likelihood of finding a map amongst a major fraction of A's children with B under C also decreases. \square

We further extend the above theorem to span sub-graph under A and B under the context of the sub-graph under C.

Theorem 2. If the bayes factor for a ontology concept A, mapped to another ontology concept B under the context of C is less than a given threshold ϵ , then the likelihood of finding a map amongst a major fraction of A's sub-graph, with a major fraction of sub-graph under B under the context of a considerable fraction of C's sub-graph also decreases.

Proof. By switching the terms A and B, in Theorem 1 we can conclude the less likelihood of finding functional links amongst children of A and B. We further include C's children, using similar logic as in Theorem 1. Since, the annotations are recursively propagated up the hierarchy in the ontology, we can apply Theorem 1 recursively to children of the concepts A, B and C, thus including the sub-graphs under them in the purview. \square

We as stated above do not give any theoretical bounds on the fraction or the likelihood of our matches, but experimentally analyze the effects of the threshold value to the running time and the amount of false negatives. The false negatives are a result of pruning the whole sub-graph under concepts. We varied our analyses over ontologies of different sizes, and observed the effects of the threshold ϵ , on the amount of pruning and the false negatives generated.

We show below the exponential reduction in running time for inferring functional links as the minimum threshold for pruning ϵ , increases.

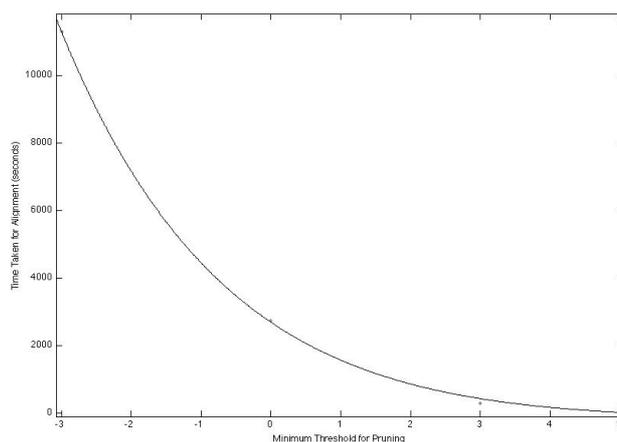


Figure 3. Graph depicting exponential reduction in running time as the minimum threshold for pruning ϵ , increases.

We also show below, the linear increase in the amount of false negatives, if we prune the full sub-graph.

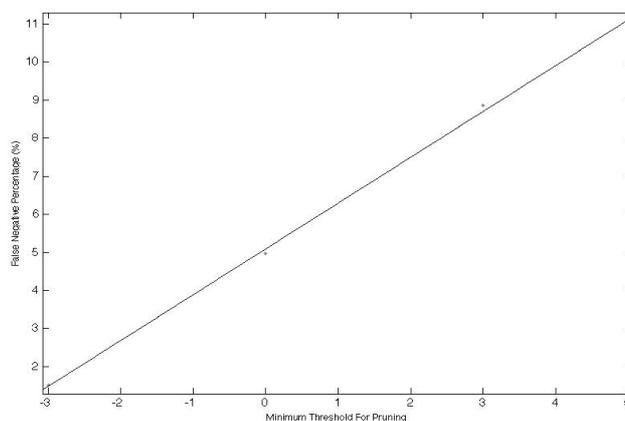


Figure 4. Graph depicting linear degradation in the amount of inferred links as the minimum threshold for pruning ϵ , increases.

We implement the heuristic pruning by performing a depth first search of the ontologies but bounding the search by the given threshold. This allows us to compute functional links with much greater efficiency with a trade-off in loss of some links.

2.4. Mapping Composition

To leverage on the inferred links, we further apply mapping composition to help us link ontologies that are not directly mapped, but have links to a common ontology under the same context. Though the algorithm is not provably composable, we apply this heuristic method to help us obtain useful connections between concepts.

To compose existing mappings, we use transitivity as the operation, that is if a link exists connecting a concept A to B, and another connecting a concept B to D under the same context C, we also map concept A to D under the context C. To validate our composed mappings, we compared them with the links found by directly mapping the ontologies.

For undirected links we found that the composed links agreed for about 45% on an average with the directly inferred links. Including directionality in the links, and then composing links with equivalent directions further increased the agreement rate to about 80%. Due to large number of ontologies, this method allows us to deduce new links by directly inferring from existing links, rather than going through the expensive process of directly mapping these ontologies.

3. Results

We obtain in all about 200 ontologies from the National Center for Biomedical Ontology's Bioportal interface. For caching sufficient statistics we obtain the dictionary of all available ontology concepts (4,153,358 terms) for searching in the corpora. We further create our b-tree index on the corpus containing the following:

1. Adverse Event Reporting System [13] database containing about 774,606 articles.
2. Array Express [14] containing 9281 articles.
3. BioSiteMaps [15] data containing 1013 articles.
4. caNanoLab [16] data containing 444 articles.
5. Conserved Domain Databases [17] containing 34,735 articles.
6. Clinical Trials [18] database containing 75,828 articles.
7. Drug Bank [19] containing 4774 articles.
8. Database of Phenotypes and Genotypes [20] having 184 articles.
9. Gene Expression Omnibus [21] containing 15,968 articles.
10. Stanford Microarray Database [22] containing 16,148 articles.
11. Published articles in PubMed [23] containing about 100,000 articles.

We then apply our Lexical Mapping algorithm, over the heuristic pruning technique to integrate Gene Ontology (containing 24,987 concepts) to all available ontologies in Bioportal under the context of the Human Disease Ontology (containing 12,033 concepts). The threshold for a significant link was set to be with a bayes factor greater than twenty ($p < 0.01$), while the threshold for pruning was set to be with a bayes factor less than zero. We further augment these links by composing the maps obtained, with Gene Ontology as the pivot. This allows us to compute mappings between any given ontologies with Human Disease as the context.

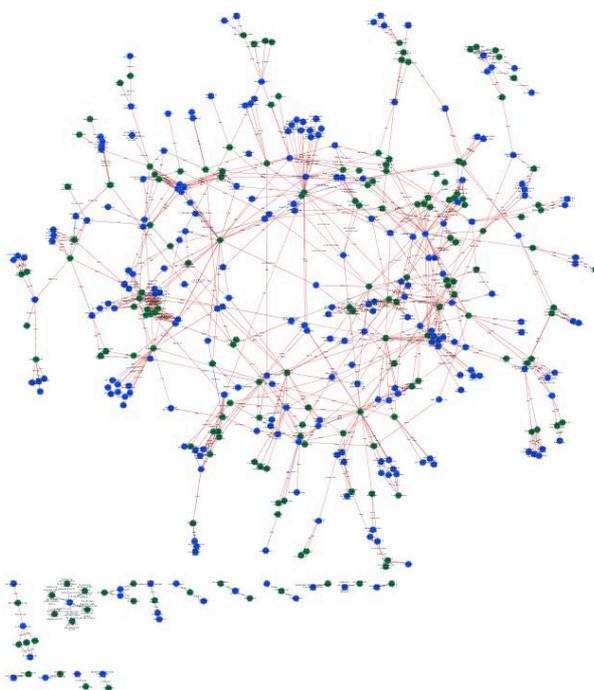


Figure 5. Mapping network showing links between Gene Ontology (blue circles) and Minimal Anatomical terminology (green circles) under Human Disease.

To validate the soundness of our inferred mappings, we take a random sampling of about 100 high information content links [24], having a significantly higher bayes factor. We then use published literature and a domain expert to biologically validate these links. The precision number for the algorithm using this approach was found to be about 0.76.

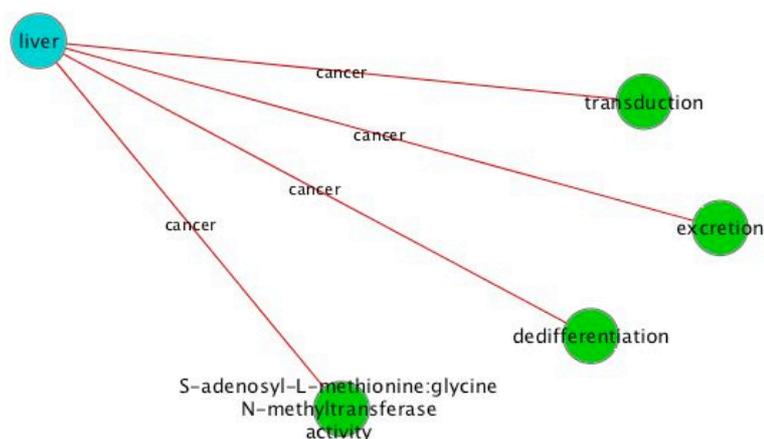


Figure 6. A portion of network showing translationalization of Gene Ontology to Anatomy under Human Disease context.

To validate the completeness of our mappings, we take a random sampling of about 100 high information content triplets of nodes. We then use published literature and a domain expert, predict links amongst these concepts. These predicted links are then matched against the ones

predicted by the algorithm to get the recall numbers. The recall number for the algorithm was found to be about 0.88. This corresponds to a f-measure of about 0.81.

4. Discussion and Conclusion

Our framework and algorithms combine disparate sources of data for discovery of relationships between ontologies. Unlike prior work, our approach tries to find context-specific functional links between ontologies, which is not possible if only semantically relevant links were considered.

By developing the novel, Lexical Mapping algorithm we identified links across ontologies, which can be used for guided expansion of various biomedical experiments. We then augmented this algorithm with heuristic approaches, for scaling up to massive data sizes with marginal loss in functional quality of links. We also introduced composition heuristics helping us to infer links between ontologies related to a different common ontology.

We further validated the utility of our algorithm, by manual verification using a domain expert, increasing confidence in our methodological approach. Our work provides a new approach for translationalizing diverse functional spaces in biomedical domain, making this huge space of knowledge amenable to researchers.

5. Acknowledgments

This work was supported in part by the National Library of Medicine (NLM/NIH) under grants 1K99LM009826 and 5T15LM007092 and by the National Human Genome Research Institute (NHGRI/NIH) under grants 2P41HG02273, 1R01HG003354, and 1R01HG004836. The authors are grateful to the anonymous reviewers for their helpful suggestions.

References

1. F. Moerchen, D. Fradkin, M. DeJori, B. Wachmann, Emerging Trend Prediction in Biomedical Literature, *AMIA Annu. Symp. Proc.*, 485 (2008).
2. M.K. Kerr, G.A. Churchill, Experimental design for gene expression microarrays, *Biostat.* 2(2): 183-201. (2001).
3. A.C. Syvänen, Toward genome-wide SNP genotyping, *Nature Gen.* 37: S5-S10. (2005).
4. B. Smith, *Ontology (Science)*, *Nature Precedings*, 2008.
5. D.L. Rubin, S.E. Lewis, C.J. Mungall, S. Misra, M. Westerfield, M. Ashburner, I. Sim, C.G. Chute, H. Solbrig, M.A. Storey, National Center for Biomedical Ontology: advancing biomedicine through structured organization of scientific knowledge, *Omics* 10(2): 185-198. (2006).
6. L. Jensen, P. Bork, *Ontologies in Quantitative Biology: A Basis for Comparison, Integration, and Discovery*, *PLoS Biol.*, (2010).
7. N. Silva, J. Rocha, Complex semantic web ontology mapping, *Web Intel. and Agent Systems* 1(3-4): 235-248. (2003).
8. A. Liddle, Information criteria for astrophysical model selection, *Monthly Notices of the Royal Astro. Soc.*, 377 (2007).
9. N. Noy, N. Shah, P. Whetzel, B. Dai, M. Dorf, N. Griffith, C. Jonquet, D. Rubin, M. Storey, C. Chute, M. Musen, *Bioportal: ontologies and integrated data resources at the click of a mouse* *Nuc. Acid. Res.*, 37 (2009).

10. C. Jonuet, N. Shah, C. Youn, M. Musen, C. Callendar, M. Storey, NCBO Annotator: Semantic Annotation of Biomedical Data, ISWC (2009).
11. M. Dai, An Efficient Solution for Mapping Free Text to Ontology Terms, AMIA Summit on Transl. Bioinfo., (2008).
12. E. Hatcher, O. Gospodnetic, Lucene in Action, JavaOne Conference, (2004).
13. S.D. Ross, M.W. Reynolds, Use of the FDA spontaneous adverse event reporting system (SAERS), or why your MedWatch reports really do matter, Journal of Clinical Oncology, 2004 ASCO Annual Meeting Proceedings (Post-Meeting Edition) 22(14S). (2004).
14. A. Brazma, H. Parkinson, U. Sarkans, M. Shojatalab, J. Vilo, N. Abeygunawardena, E. Holloway, M. Kapushesky, P. Kemmeren, G.G. Lara, A. Oezcimen, P. Rocca-Serra, S.A. Sansone, ArrayExpress—a public repository for microarray gene expression data at the EBI, Nucleic Acids Res. 31(1):68-71. (2003).
15. L. Marenco, R. Wang, G.M. Shepherd, P.L. Miller, The NIF DISCO Framework: Facilitating Automated Integration of Neuroscience Content on the Web, Neuroinform 8:101-112. (2010).
16. V. Maojo, F. Martin-Sanchez, C. Kulikowski, A. Rodriguez-Paton, M. Fritts, Nanoinformatics and DNA-Based Computing: Catalyzing Nanomedicine, Ped. Research 67(5): 481-489. (2010).
17. A. Marchler-Bauer, J.B. Anderson, F. Chitsaz, M.K. Derbyshire, C. DeWeese-Scott, J.H. Fong, L.Y. Geer, R.C. Geer, N.R. Gonzales, M. Gwadz, S. He, D.I. Hurwitz, J.D. Jackson, Z. Ke, C.J. Lanczycki, C.A. Liebert, C. Liu, F. Lu, S. Lu, G.H. Marchler, M. Mullokandov, J.S. Song, A. Tasneem, N. Thanki, R.A. Yamashita, D. Zhang, N. Zhang, S.H. Bryant, CDD: specific functional annotation with the Conserved Domain Database, Nucleic Acids Res. 37: D205-10. (2009).
18. M. Mi, Clinical Trials Database: Linking Patients to Medical Research <http://clinicaltrials.gov>, Journal of Consumer Health On the Internet, 9(3): 59-67. (2005).
19. D.S. Wishart, C. Knox, A.C. Guo, D. Cheng, S. Shrivastava, D. Tzur, B. Gautam, M. Hassanali, DrugBank: a knowledgebase for drugs, drug actions and drug targets, Nucleic Acids Res. 36(Database issue):D901-6. (2008).
20. M.D. Mailman, M. Feolo, Y. Jin, M. Kimura, K. Tryka, R. Bagoutdinov, L. Hao, A. Kiang, J. Paschall, L. Phan, N. Popova, S. Pretel, L. Ziyabari, Y. Shao, Z.Y. Wang, K. Sirotkin, M. Ward, M. Kholodov, K. Zbicz, J. Beck, M. Kimelman, S. Shevelev, D. Preuss, E. Yaschenko, A. Graeff, J. Ostell, S.T. Sherry, The NCBI dbGaP database of genotypes and phenotypes, Nat. Genet. 39(10): 1181-1186. (2007).
21. T. Barrett, D.B. Troup, S.E. Wilhite, P. Ledoux, D. Rudnev, C. Evangelista, I.F. Kim, A. Soboleva, M. Tomashevsky, K.A. Marshall, K.H. Phillipy, P.M. Sherman, R.N. Muerter, R. Edgar, NCBI GEO: archive for high-throughput functional genomic data, Nucleic Acids Res. 37(Database issue):D5-15. (2009).
22. J. Hubble, J. Demeter, H. Jin, M. Mao, M. Nitzberg, T.B. Reddy, F. Wymore, Z.K. Zachariah, G. Sherlock, C.A. Ball, Implementation of GenePattern within the Stanford Microarray Database. Nucleic Acids Res. 37(Database Issue):D898-901. (2009).
23. C. Plake, T. Schiemann, M. Pankalla, J. Hakenberg, U. Leser, ALIBABA: PubMed as a graph, Bioinformatics 22(19):2444-2445. (2006).

24. G. Alterovitz, M. Xiang, D.P. Hill, J. Lomax, J. Liu, M. Cherkassky, J. Dreyfuss, C. Mungall, M.A. Harris, M.E. Dolan, J.A. Blake, M.F. Ramoni, Ontology engineering, *Nat. Biotech.* 28:128-130. (2010).