# The coevolution of ontologies and knowledge-based analytics in bioinformatics

Robert Hoehndorf[1]

[1]*Computational Bioscience Research Center, Computer, Electrical and Mathematical Sciences & Engineering Division, King Abdullah University of Science and Technology, 4700 KAUST, Thuwal 23955, Saudi Arabia*

## Abstract

I discuss the coevolution of bio-ontologies and analytical bioinformatics methods in response to the evolving landscape of life sciences. I focus on the role of ontologies, in particular the Gene Ontology, in capturing and describing biological knowledge, and the challenges and developments in ontology-based bioinformatics, particularly in light of new computational methods and machine learning. The main theme is the bidirectional influence between how ontologies and bioinformatics methods evolved together, and how ontologies have shaped advancements in the analysis, representation, and understanding of biological data by providing a unifying layer of knowledge.

## Keywords

bio-ontology, knowledge-based analytics, Artificial Intelligence

Developing high quality ontologies is expensive, and, like most infrastructure components of the life sciences, ontologies have evolved in response to specific needs and requirements of the biomedical community. At the same time, new tools utilizing ontologies emerged to enable or improve analysis of biological data. In my talk, I will explore how bio-ontologies have evolved in response to a changing bioinformatics environment and how bioinformatics tools and methods evolved in response to changing ontologies; my main aim will be to characterize the current changes in bioinformatics through large-scale application of machine learning, and how ontologies have to change to accommodate these changes.

The Gene Ontology (GO), the first bio-ontology that was and still is widely used, emerged as a consequence of breakthroughs in gene and genome sequencing and the resulting understanding of how many genes are conserved in different organisms [1]. This novel understanding, combined with the rapid change of knowledge in the field of molecular biology, necessitated the development of the GO, to keep track of the changing knowledge in the field and simultaneously provide a means to describe our knowledge of gene and protein functions. Using the GO for describing protein functions solved many challenges. A form of deductive inference ("true path rule") allowed capturing the most specific information about a protein as possible while still allowing inference of more general information, and use of a taxonomy allowed knowledge to evolve by gradually adding more specific functions to a protein without invalidating previous assertions.

Today, some of the most exciting developments (and

challenges) in bio-ontologies still occur in fields where novel experimental techniques are leading to a radical change of our understanding of biological phenomena. For example, recently, our understanding of cell types has changed drastically, resulting from single cell sequencing technologies and the resulting detailed information available about cell types and their relations; ontologies of cell types had to change accordingly [2], and cell ontologies are now one of the most active areas of bio-ontology development (as evidenced, for example, by the regular CELLS workshop co-lated with the International Conference on Biomedical Ontologies).

Yet, what the early development of the GO (and similar ontologies) has shown is that the development and evolution of ontologies in life sciences is not a one-way road and only determined by changes in experimental techniques; rather, the availability of ontologies has also led to novel computational analysis methods, and ontologies will change in response to the emergence of novel methods. Two methods are particularly noteworthy here, ontology enrichment analysis and semantic similarity measures. Both techniques are some of the most widely used computational analysis methods involving ontologies. An ontology enrichment analysis uses an ontology together with its annotations in order to determine whether there is a function that is statistically enriched in a set of genes or gene products [3, 4, 5]. Ontology-based semantic similarity measures utilize the knowledge contained in ontologies (in particular within the formal axioms) to define measures of similarities between ontology classes, sets of classes, instances of classes, or entities annotated with (sets of) classes [6]. Semantic similarity was first used to query for and retrieve "semantically" related proteins [7], and later extended to find other entities with some association using a "guilt by association"

approach.

My key take away message from these methods is that bioinformatics has developed a set of computational methods that crucially relied on ontologies providing accurate results. Both enrichment analysis and semantic similarity require that inferences in ontologies, in particular inferences about annotated genes or gene products (the "true path rule" in GO and more elaborate versions of this rule), are *accurate* (accurate in the sense that they are biologically correct and experimentally verifiable). Early ontologies did not always produce accurate inferences [8, 9, 10, 11], and finding these incorrect inferences has, arguably, led to one of the most active periods for ontology development and quality improvement, where the community applied and developed methods inspired, among others, by philosophy [12], linguistics [13], and logics [14].

With further improvement in experimental methods, in particular the emergence of high throughput sequencing methods, the demands on ontologies rose further, both in terms of their accuracy as well as in their detail and discriminatory power. Ontologies now had to cope with Big Data, and manually building ontologies would no longer scale in many domains. In this time, ontology design patterns [15], upper ontologies [16, 17], more and elaborate ontology design principles and community standards allowed ontologies to "scale up" both to capturing Big Data and more detailed nuances in biological phenomena. The new problem arose that our tools (reasoners and ontology editors such as Protege) no longer scaled to the new size and complexity of ontologies. The solution was to switch to different tools like Elk [18], and apply modularization techniques such as MIREOT [19]; while these work in solving the problem of scalability to Big Data, they have also hidden (and lost) some information; automated reasoners such as Elk only consider a tiny subset of the language we use to formalize ontologies, and modularization techniques can hide inconsistencies and therefore allow inconsistencies to increase [20].

As a result, a switch took place within the bio-ontologies community and the focus was no longer only on "ontologies" as formal artifacts capturing domain knowledge accurately, but rather on constructing "knowledge graphs" in which the focus is on linking information in some (vaguely) meaningful manner. The tendency to focus more on "knowledge graphs" instead of ontologies was by no means universal but certainly noticeable and still ongoing today. The move was motivated by the desire to focus on "relatedness" instead of precision, and find ways to integrate (i.e., link) large amounts of resources, in particular in the biomedical domain; the resources that are linked were often not ontologies but (medical) terminologies, so that ontological precision may have been an obstacle to successful integration.

At the same time, and further motivating the focus on knowledge graphs instead of ontologies, novel knowledge graph analytics approaches emerged, in particular machine learning methods that would operate directly on graphs or knowledge graphs [21, 22], and graph neural networks that can exploit the knowledge graphs for various tasks [23]. In particular, knowledge graph embedding methods have been adopted widely within the bioinformatics community to exploit information in knowledge graphs for predictive or analytical tasks. Several knowledge graph embedding methods have been developed [21], but some of the most popular are based on the principle that, if the fact $r(a, b)$ is in the knowledge graph, then $\vec{a} + \vec{r} \approx \vec{b}$ (where $\vec{a}$ etc. are the "embedding" vectors of some dimension that "represent" $a$, $r$, and $b$ in a distributed manner) [24]. The advantage of these embedding methods is their interpretability, simplicity, and almost universal applicability.

The role of ontologies in graph-based machine learning methods (such as knowledge graph embeddings, or graph neural networks) is to provide a source of nodes, and the formal axioms in the ontologies provides a source of relatedness (edges) that make up the resulting graph [25]. Yet, many aspects that have been considered crucial in developing ontologies are lost, specifically all benefits arising from semantics, both logical and ontological [26]: the ability for complex queries; ensured consistency; and deductive inference. In particular deductive inference (which is required both for complex queries and determining consistency) is crucial for exploring the knowledge ontologies contain beyond what has been explicitly asserted, but this ability for deductive inference is largely lost in graph-based methods.

Before ontologies (considered here as artifacts which explicitly and formally specify a conceptualization of a domain using a logic-based language) can become relevant in machine learning in bioinformatics, methods that can utilize the semantics of ontologies need to first be developed, because very few such methods exist in the field of AI; and it is even more of a challenge to tune such methods to the specific peculiarities of bio-ontologies which have distinct properties when compared to ontologies used in other domains, in particular computer science.

Some new methods emerged over the past years that apply machine learning methods to bio-ontologies. While some of these methods are simple extensions of learning from graph-structured data or learning from text, more recent approaches aim to explicitly address the missing formal semantics in machine learning models. These neuro-symbolic methods can produce deductive inferences directly, either by implementing a deduction system using neural approaches or by generating model structures using neural approaches. Establishing this correspondence between classical semantics and neural

networks enables novel applications and demands on ontologies, but also opens novel opportunities, both for bioinformatics and AI. In bioinformatics, these methods allow machine learning to utilize the vast and rich knowledge contained in bio-ontologies thereby endowing the machine learning models with domain knowledge (and the ability to explore the knowledge more deeply than would be possible using only knowledge graphs), which can be used to provide access to the results of over a hundred years of experiments that are now contained in ontologies and knowledge bases. One of the most obvious areas of application are rare diseases where only little training data will ever be available. For AI, bio-ontologies provide a vast und largely underused resource of knowledge with direct implications for health, the environment, and well-being.

# References

[1] M. Ashburner, C. A. Ball, J. A. Blake, D. Botstein, H. Butler, M. J. Cherry, A. P. Davis, K. Dolinski, S. S. Dwight, J. T. Eppig, M. A. Harris, D. P. Hill, L. I. Tarver, A. Kasarskis, S. Lewis, J. C. Matese, J. E. Richardson, M. Ringwald, G. M. Rubin, G. Sherlock, Gene ontology: tool for the unification of biology, Nature Genetics 25 (2000) 25–29. URL: http://dx.doi.org/10.1038/75556. doi:10.1038/75556.

[2] D. Osumi-Sutherland, C. Xu, M. Keays, A. P. Levine, P. V. Kharchenko, A. Regev, E. Lein, S. A. Teichmann, Cell type ontologies of the human cell atlas, Nature Cell Biology 23 (2021) 1129–1135. URL: https://doi.org/10.1038/s41556-021-00787-7. doi:10.1038/s41556-021-00787-7.

[3] S. W. Doniger, N. Salomonis, K. D. Dahlquist, K. Vranizan, S. C. Lawlor, B. R. Conklin, MAPPFinder: using Gene Ontology and GenMAPP to create a global gene-expression profile from microarray data, Genome Biology 4 (2003) R7. URL: https://doi.org/10.1186/gb-2003-4-1-r7. doi:10.1186/gb-2003-4-1-r7.

[4] A. Subramanian, P. Tamayo, V. K. Mootha, S. Mukherjee, B. L. Ebert, M. A. Gillette, A. Paulovich, S. L. Pomeroy, T. R. Golub, E. S. Lander, J. P. Mesirov, Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles, Proceedings of the National Academy of Sciences of the United States of America 102 (2005) 15545–15550. URL: http://www.pnas.org/content/102/43/15545.abstract. doi:10.1073/pnas.0506580102.

[5] M. D. Robinson, J. Grigull, N. Mohammad, T. R. Hughes, FunSpec: a web-based cluster interpreter for yeast, BMC Bioinformatics 3 (2002) 35. URL: https://doi.org/10.1186/1471-2105-3-35. doi:10.1186/1471-2105-3-35.

[6] S. Harispe, S. Ranwez, S. Janaqi, J. Montmain, The semantic measures library and toolkit: fast computation of semantic similarity and relatedness using biomedical ontologies, Bioinformatics 30 (2014) 740–742. URL: http://bioinformatics.oxfordjournals.org/content/30/5/740.abstract. doi:10.1093/bioinformatics/btt581.

[7] P. W. Lord, R. D. Stevens, A. Brass, C. A. Goble, Investigating semantic similarity measures across the gene ontology: the relationship between sequence and annotation, Bioinformatics 19 (2003) 1275–1283. URL: http://bioinformatics.oxfordjournals.org/content/19/10/1275.abstract. doi:10.1093/bioinformatics/btg153.

[8] B. Smith, J. Williams, S. Schulze-Kremer, The ontology of the gene ontology., AMIA ... Annual Symposium proceedings / AMIA Symposium. AMIA Symposium (2003) 609–613. URL: http://view.ncbi.nlm.nih.gov/pubmed/14728245.

[9] B. Smith, C. Rosse, The role of foundational relations in the alignment of biomedical ontologies., Medinfo 11 (2004) 444–448.

[10] B. Smith, Against fantology, in: M. E. Reicher, J. C. Marek (Eds.), Experience and Analysis. Proceedings of the 27th International Wittgenstein Symposium., volume 6, 2005, pp. 153–170. URL: http://dx.doi.org/10.1186/gb-2004-6-1-r7. doi:http://dx.doi.org/10.1186/gb-2004-6-1-r7.

[11] W. Ceusters, P. Elkin, B. Smith, Referent tracking: The problem of negative findings, Stud Health Technol Inform (2006).

[12] B. Smith, W. Ceusters, Ontological realism: A methodology for coordinated evolution of scientific ontologies, Appl. Ontol. 5 (2010) 139–188.

[13] M. Bada, L. Hunter, Enrichment of OBO ontologies, Journal of Biomedical Informatics 40 (2007) 300–315. URL: https://doi.org/10.1016/j.jbi.2006.07.003. doi:10.1016/j.jbi.2006.07.003.

[14] R. Hoehndorf, F. Loebe, J. Kelso, H. Herre, Representing default knowledge in biomedical ontologies: application to the integration of anatomy and phenotype ontologies, BMC Bioinform. 8 (2007). URL: https://doi.org/10.1186/1471-2105-8-377. doi:10.1186/1471-2105-8-377.

[15] D. Osumi-Sutherland, M. Courtot, J. P. Balhoff, C. Mungall, Dead simple OWL design patterns 8 (2017). URL: https://doi.org/10.1186/s13326-017-0126-0. doi:10.1186/s13326-017-0126-0.

[16] B. Smith, W. Ceusters, B. Klagges, J. Köhler, A. Kumar, J. Lomax, C. Mungall, F. Neuhaus, A. L. Rector, C. Rosse, Relations in biomedical ontologies., Genome Biol 6 (2005) R46. URL: http://dx.doi.org/10.

1186/gb-2005-6-5-r46. doi:http://dx.doi.org/10.1186/gb-2005-6-5-r46.

[17] B. Smith, M. Ashburner, C. Rosse, J. Bard, W. Bug, W. Ceusters, L. J. Goldberg, K. Eilbeck, A. Ireland, C. J. Mungall, N. Leontis, P. R. Serra, A. Ruttenberg, S. A. Sansone, R. H. Scheuermann, N. Shah, P. L. Whetzel, S. Lewis, The OBO Foundry: coordinated evolution of ontologies to support biomedical data integration, Nat Biotech 25 (2007) 1251–1255.

[18] Y. Kazakov, M. Krötzsch, F. Simancik, The incredible elk, Journal of Automated Reasoning 53 (2014) 1–61. URL: http://dx.doi.org/10.1007/s10817-013-9296-3. doi:10.1007/s10817-013-9296-3.

[19] M. Courtot, N. Juty, C. Knüpfer, D. Waltemath, A. Zhukova, A. Dräger, M. Dumontier, A. Finney, M. Golebiewski, J. Hastings, S. Hoops, S. Keating, D. B. Kell, S. Kerrien, J. Lawson, A. Lister, J. Lu, R. Machne, P. Mendes, M. Pocock, N. Rodriguez, A. Villeger, D. J. Wilkinson, S. Wimalaratne, C. Laibe, M. Hucka, N. Le Novère, Controlled vocabularies and semantics in systems biology., Molecular systems biology 7 (2011). URL: http://dx.doi.org/10.1038/msb.2011.77. doi:10.1038/msb.2011.77.

[20] L. T. Slater, G. V. Gkoutos, R. Hoehndorf, Towards semantic interoperability: finding and repairing hidden contradictions in biomedical ontologies, BMC Medical Informatics and Decision Making 20 (2020). URL: https://doi.org/10.1186/s12911-020-01336-2. doi:10.1186/s12911-020-01336-2.

[21] Q. Wang, Z. Mao, B. Wang, L. Guo, Knowledge graph embedding: A survey of approaches and applications, IEEE Transactions on Knowledge and Data Engineering 29 (2017) 2724–2743. doi:10.1109/TKDE.2017.2754499.

[22] M. Ali, C. T. Hoyt, D. Domingo-Fernández, J. Lehmann, H. Jabeen, BioKEEN: a library for learning and evaluating biological knowledge graph embeddings, Bioinformatics 35 (2019) 3538–3540. URL: https://doi.org/10.1093/bioinformatics/btz117. doi:10.1093/bioinformatics/btz117.

[23] X.-M. Zhang, L. Liang, L. Liu, M.-J. Tang, Graph neural networks and their current applications in bioinformatics, Frontiers in Genetics 12 (2021). URL: https://doi.org/10.3389/fgene.2021.690049. doi:10.3389/fgene.2021.690049.

[24] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, O. Yakhnenko, Translating embeddings for modeling multi-relational data, Advances in neural information processing systems 26 (2013).

[25] J. Chen, P. Hu, E. Jimenez-Ruiz, O. M. Holter, D. Antonyrajah, I. Horrocks, OWL2Vec*: embedding of OWL ontologies, Machine Learning (2021). URL: https://doi.org/10.1007/s10994-021-05997-6. doi:10.1007/s10994-021-05997-6.

[26] F. Loebe, H. Herre, Formal semantics and ontologies - towards an ontological account of formal semantics, in: C. Eschenbach, M. Grüninger (Eds.), Formal Ontology in Information Systems, Proceedings of the Fifth International Conference, FOIS 2008, Saarbrücken, Germany, October 31st - November 3rd, 2008, volume 183 of *Frontiers in Artificial Intelligence and Applications*, IOS Press, 2008, pp. 49–62. URL: https://doi.org/10.3233/978-1-58603-923-3-49. doi:10.3233/978-1-58603-923-3-49.