

Ontology Dynamics in a Data Life Cycle: Challenges and Recommendations from a Geoscience Perspective

Xiaogang Ma*, Peter Fox, Eric Rozell, Patrick West, Stephan Zednik

Tetherless World Constellation, School of Science, Rensselaer Polytechnic Institute, 110 8th Street, Troy, NY 12180, USA

ABSTRACT: Ontologies are increasingly deployed as a computer-accessible representation of key semantics in various parts of a data life cycle and, thus, ontology dynamics may pose challenges to data management and re-use. By using examples in the field of geosciences, we analyze challenges raised by ontology dynamics, such as heavy reworking of data, semantic heterogeneity among data providers and users, and error propagation in cross-discipline data discovery and re-use. We also make recommendations to address these challenges: (1) communities of practice on ontologies to reduce inconsistency and duplicated efforts; (2) use ontologies in the procedure of data collection and make them accessible to data users; and (3) seek methods to speed up the reworking of data in a Semantic Web context.

KEY WORDS: semantic web, knowledge evolution, data transformation, geoscience.

1 INTRODUCTION

A short article (Mascarelli, 2009) published in Nature in 2009 reveals the geoscience community's concern of the potential heavy reworking of data due to the changes to the definition of a geological time term (i.e., Quaternary) made by the International Commission on Stratigraphy (ICS). The article reported that, in the 1980s, the US Geological Survey had reworked all of its maps and terminology because the definition of another geological time term (i.e., Pleistocene) was changed. Following the changes made to Quaternary, the data rework might to be performed again. Extending our view from the nomenclature changes and data rework discussed in that article to the research on ontology dynamics, we can see a topic of interest is to consider ontology dynamics within a data life cycle.

Ontology in computer science is defined as a shared conceptualization of domain knowledge (Guarino, 1997; Gruber, 1995). The evolution of knowledge within a domain is continuous, which brings dynamics to ontologies of that domain. Ontology dynamics is understood as an umbrella field that covers research topics like ontology versioning, inconsistency handling, ontology evolution and integration, and change propagation, etc. (Flouris et al., 2009). A data life cycle represents the whole procedure of data management (an example is shown in Fig. 1). Ontology can be, but is not always, a primary part in the first step (i.e., concept articulation, definition, etc.) of the data life cycle. Any changes in the first step may require modifications in knowledge encodings in the following steps

to make the concepts consistent across the cycle.

Turning our view back to the field of geosciences, we can see that the relationship between ontology dynamics and a data life cycle has not obtained enough attention yet (Geo-Data Informatics Workshop Committee, 2011). In recent years, geoscience data are increasingly shared online. However, metadata like data provenance, data structure, and ontology use are often absent or incomplete (Ma et al., 2011a). With the above-mentioned examples in the Nature magazine, we can see that the ontologies of a geoscience data provider may be changed and the data updated, but a general user may not know of these changes, let alone the implications. As a result, the user may understand the data in a different meaning that the data provider intends. Ontology dynamics thus raises difficulties to a data life cycle. The problem may be even worse if a researcher is using data provided by (for example) a web service that uses an inconsistent knowledge encoding and “hides” that fact from the researcher.

The purpose of this paper is to discuss challenges that ontology dynamics may pose to the data life cycle in geosciences, and to make recommendations on how to address these challenges. To support the discussion, examples of knowledge evolution and ontology dynamics surrounding the topic of geological time scale will be used in the following sections. Although these examples are in the field of geosciences, they share mutual interests with studies on knowledge evolution, ontology dynamics and other semantic technologies in the field of computer science. Addressing the challenges discussed here calls further efforts from both the geoscience and the computer science communities. In Section 2 we present examples of ontology dynamics surrounding the topic of geological time scale. In Section 3 we discuss the challenges derived from ontology dynamics, and make recommendations to address them. Section 4 wraps up the paper and draws conclusions.

*Corresponding author: max7@rpi.edu

© China University of Geosciences and Springer-Verlag Berlin Heidelberg 2014

Manuscript received September 19, 2013.

Manuscript accepted March 3, 2014.

2 AN ONTOLOGY SPECTRUM OF GEOLOGICAL TIME SCALE

The geological time scale is a chronological framework for studying the history of the Earth. The ICS works on uniting local and regional classifications of geological time scales in order to establish a standard time scale for global correlations. ICS advances its work by reaching agreements on nomenclature and a hierarchy of temporal intervals and boundaries defined by Global Boundary Stratotype Sections and Points (GSSPs) (International Commission on Stratigraphy, 2012). The International Stratigraphic Chart, an output from these works, is used in geological works across the world. Following the progress of ICS, the chart has been updated frequently (Fig. 2).

The work of ICS reflects knowledge evolution of the geological time scale due to the progress of studies in stratigraphy,

and there are also examples that reflect ontology dynamics surrounding the same topic: glossaries (Bates and Jackson, 1995), taxonomies (British Geological Survey, 2012), thesauri (CCOP and CIFEG, 2006) and controlled vocabularies (Ma et al., 2010) that include terms of the geological time scale. There are also conceptual schemas expressed using the Unified Modeling Language (UML) (Cox and Richard, 2005; NADM Steering Committee, 2004) and the Geography Markup Language (GML) (Sen and Duffy, 2005). Also there are ontologies of the geological time scale that use Semantic Web languages such as the Simple Knowledge Organization System (SKOS) (Ma et al., 2011b) and the Web Ontology Language (OWL) (Ma and Fox, 2013; Cox, 2011). These works constitute a series of ontologies with different levels of semantic richness, named by computer scientists as an ontology spectrum (Fig. 3) (McGuinness, 2003; Obrst, 2003).

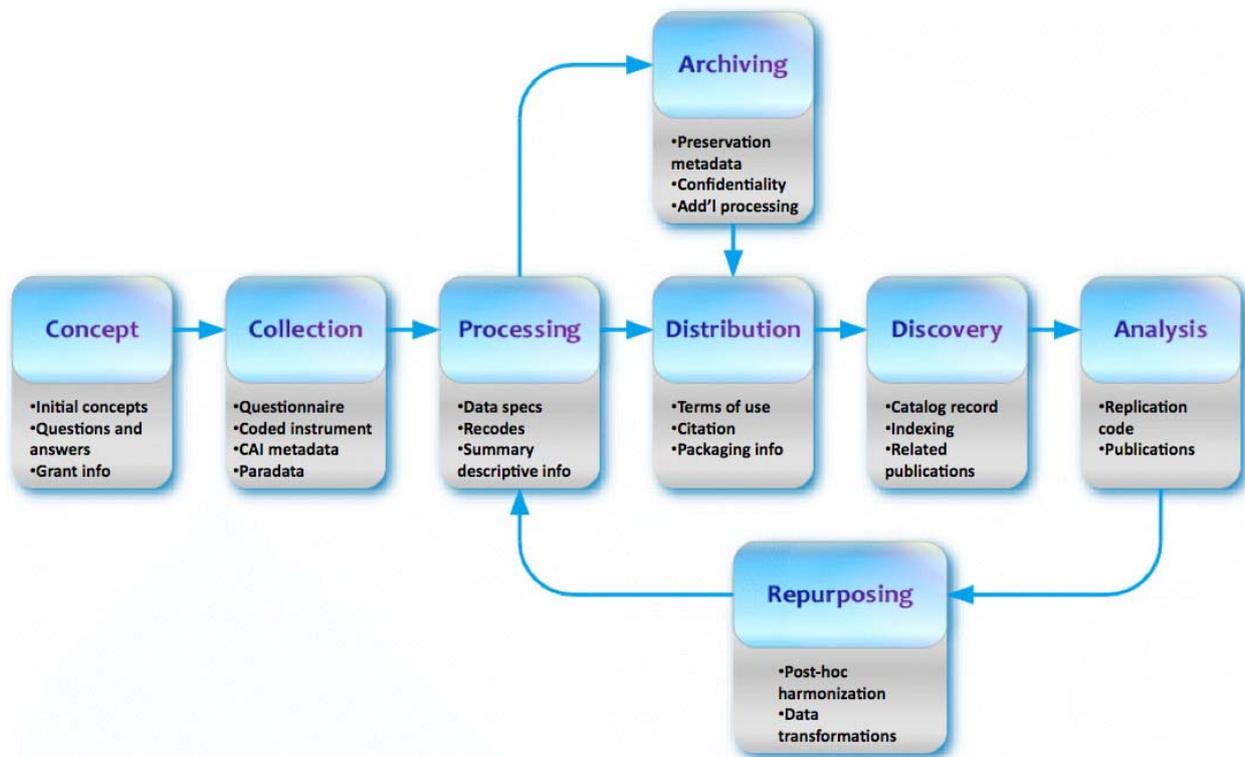


Figure 1. Data life cycle (from Spencer, 2012).

Most recently, researchers have built ontologies in Semantic Web languages for different versions of the International Stratigraphic Chart (Cox, 2012). Listings 1 and 2 (see Appendix) are source codes of a same concept from two ontologies in Cox (2012). They reflect a part of the change in the International Stratigraphic Chart from version 2008 to 2009 (see Fig. 2). That is, the base of Quaternary is changed from 1.806 to 2.588 Ma (here Ma can be understood as 'million years ago'), so the GSSP of Quaternary base is also shifted from the base of Calabrian to the base of Gelasian, as described in Listings 1 and 2, respectively. Another feature of these ontologies (Cox, 2012), is that the concept definitions also provide mappings to equivalent concepts in other ontology

works such as SWEET (Raskin and Pan, 2005) and DBpedia by using `<owl:sameAs>` and `<foaf:isPrimaryTopicof>`.

Through the above-mentioned examples we can see that, for a single domain (i.e., geological time scale), there are several ontology types with different semantic richness. Even if we focus on certain types (i.e., those encoded in Semantic Web languages), there exist different versions of ontologies (i.e., they present actual examples of the ontology dynamics of interest herein). A recent review (Flouris et al., 2008) summarizes 10 topics of general ontology changes/dynamics such as: ontology mapping, morphism, evolution, debugging and versioning, etc.. Actual examples of evolution, versioning and mapping can be found in the above ontology works on the

geological time scale.

In the field of computer science, works addressing ontology dynamics have made progress in several directions, such as developing a uniform framework for managing ontology versions (Noy and Musen, 2004), inconsistency handling by resolving or reasoning (Bell et al., 2007) and query answering under evolving ontologies to avoid redefining mappings between ontologies and data sources (Kondylakis and Plexousakis, 2011), etc.. Nonetheless, there is a shortage of works that consider ontology dynamics consistently within a data life cycle. Several ontologies in the above-mentioned examples of geological time have been used in actual works to promote data, management, standardization and integration. For example, the controlled vocabulary in Ma et al. (2010) was used to standardize terminology in databases and enable semantic interoperability of data for mining projects. The NADM

model (NADM Steering Committee, 2004) was used in the US National Geological Map Database project as a reference standard (Soller and Berg, 2005). The GeoSciML (Sen and Duffy, 2005) and its corresponding vocabularies were used in the OneGeology and OneGeology-Europe projects for harmonizing distributed geological maps (Laxton et al., 2010). Despite such progress, discussion of these works, in the context of ontology dynamics and its impacts on a data life cycle, are significantly less developed and less formal in the field of geosciences, compared to those in computer science.

3 CHALLENGES AND RECOMMENDATIONS REGARDING A DATA LIFE CYCLE

We now discuss some challenges that ontology dynamics may pose to the works of a data life cycle, using examples in geosciences.

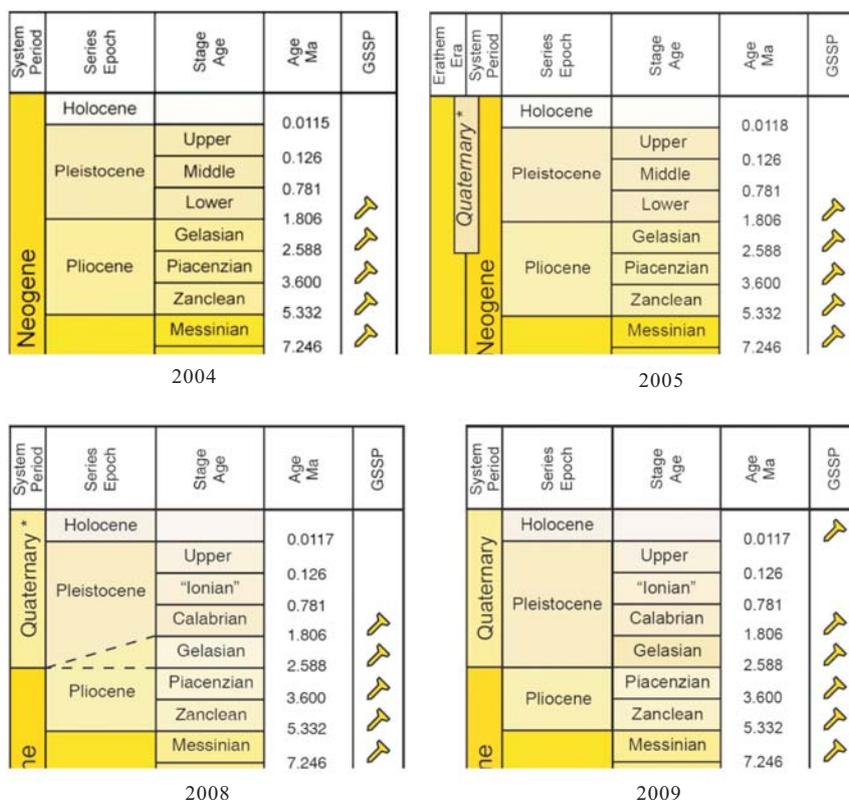


Figure 2. Definition of Quaternary and its sub-concepts in different versions of the International Stratigraphic Chart.

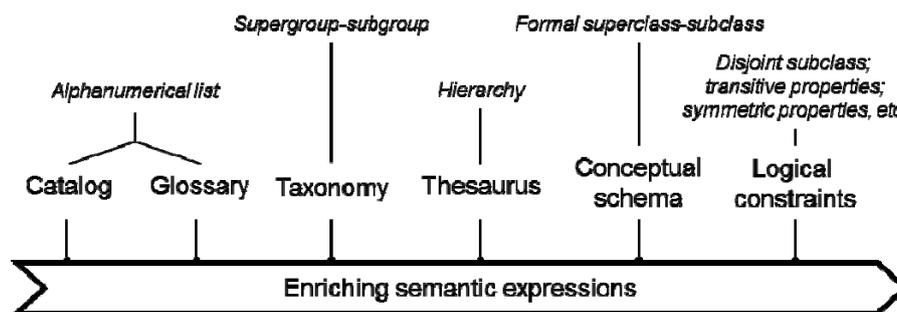


Figure 3. Adaptation of the ontology spectrum from McGuinness (2003) and Obrst (2003).

The first challenge is caused by ontology versioning. Geoscience data include both spatial and attribute records, which can be difficult to rework. For example, on a map, the surficial rock age of an area may originally be recorded as Neogene, whereas now a part of the area may be recorded as Quaternary following the changes described in Mascarelli (2009). Moreover, the polygon of the original area on the map will be split, and the part characterized as Quaternary may be merged with an adjacent polygon of Quaternary. A map may contain, depending on the scale, tens to hundreds or thousands of such polygons requiring reworking, and there may be tens to hundreds of maps for a nation or region. Because the ICS-specified concepts of geological time scale are accepted and used globally, and geological time is an essential part of basic geological maps, the data reworking in Mascarelli (2009) would be substantial and potentially prone to errors, if it is considered at a global level. Further, products derived from these types of map would thus, by definition, be 'out-of-date'.

The second challenge is of semantic mismatch caused by ontology inconsistency between data providers. Different ontologies of the same domain are often used in the data of different institutions (often in different countries). Definitions of the same concept may be diverse among those ontologies and, in turn, in the data of those data providers. Even for the data collected and archived by the same institution, there may be several versions of ontologies on a same topic, so the semantic mismatch may also exist within a single institution.

The third challenge is the semantic mismatch between data providers and data users. For example, a data provider understands Quaternary as a period from present to 1.806 Ma, while a data user understands it as from present to 2.588 Ma. The difference is 0.782 Ma, which may bring huge changes or errors to the results of the data user. Moreover, if a user's data are retrieved from multiple sources with mismatched ontologies among them, the challenge will be even more serious and often impossible to reconcile.

We characterize the fourth challenge as error propagation, arising in cross-discipline data discovery and re-use. It may be included in the third challenge, but it is worth being listed independently. To address societal challenges like climate change and environmental degradation, data from different disciplines are often assembled, whereas a data user may not know the background knowledge of each discipline. The errors described in the third challenge may happen in a step of the data processing, and the result will then be used in the next step. Finally, there will be deviations in the overall result of the, otherwise desirable, multi-disciplinary study. For example, a recent study (Allen et al., 2008) described the large, multi-source, multi-disciplinary and inconsistent datasets of the cross-border Abbotsford-Sumas aquifer between British Columbia, Canada and Washington State, USA. The data in that study vary in style, classification and nomenclature. The authors listed a few examples, such as "the geological classification for the bridge construction reports were based on the Unified Soil Classification System, which is used for engineering purposes and is based on the particle size, liquid limit and plasticity index; the drill core record descriptions were based on the Wentworth Scale; GSC geological descriptions were

based on the stratigraphy and environment of deposition, and the drillers' descriptions were based on experience or education." Such variation resulted in over 6 000 unique geological categories in the water well records of that region (Allen et al., 2008). Though not reported in that study but we can see the heterogeneous datasets may cause misconception in various subsequent works, because groundwater is not only used for domestic life, but also supports industrial, farming and agricultural activities in that area.

Some of these challenges have existed for a long term and may never be fully solved. We would like to present the following recommendations to address these challenges or to reduce their negative impacts.

The first recommendation is for communities of practice and collaborative work on ontologies. We described above that for the topic of geological time scale, there are several types of ontologies and some of them have several versions. Because in general, ontologies in geosciences are still underdeveloped, many data providers are building so-called micro-ontologies (Janowicz and Hitzler, 2012) to support their particular working tasks. While accepting that diversity (i.e., a bottom-up approach) is an engine of scientific works, we argue that collaborative (i.e., top-down) approaches can bring together people under the premise of mutual benefit in the same or similar discipline to share approaches and lessons learned, with the aims to reduce inconsistency and duplicated efforts. Moreover, there should also be some collaboration among these domain-dependent efforts and other coordinating and standards organizations, such as the Open Geospatial Consortium (OGC) and the World Wide Web Consortium (W3C), to mitigate conflicts. Consequently, inconsistency between data retrieved from different sources can be reduced. For example, the GeoSciML (Laxton et al., 2010; Sen and Duffy, 2005) and its corresponding vocabularies are results of a procedure that integrates both bottom-up and top-down approaches, implemented by researchers from geological survey organizations across the world. Those researchers participate in an organization called the Commission for the Management and Application of Geoscience Information under the International Union of Geological Sciences (CGI-IUGS) to coordinate works surrounding the GeoSciML and the vocabularies. People in CGI-IUGS also work together with OGC and W3C as well as several other associations and initiatives, such as the Federation of Earth Science Information Partners (ESIP) and the Infrastructure for Spatial Information in the European Community (INSPIRE), to promote the interoperability of the developed geoscience ontologies and vocabularies and the data underpinned by them.

The second recommendation is to use standardized terms in the procedure of data collection and provide data users the access to ontologies that give detailed annotations of those terms and their inter-relationships. Automatic and semi-automatic tools (both software and hardware) can be developed to promote the use of ontologies in data collection and archiving, thus reduce conceptual inconsistency within an institution. For example, the *FieldLog* system (Brodaric, 2004) provides a common ontology consisting of cartographic, geospatial, geological and metadata elements, and associated attributes, constraints, behavior and relations, which enables individual data-

base designs to be created to support the collection of geological data in the field. If a number of ontologies are used in the collection of data, then they can also be archived, for example, as a part of the metadata together with the data. For example, in the online geologic maps of US states (<http://mrddata.usgs.gov/geology/state>), each recoded term of rock type has a hyperlink directing to a vocabulary page with detailed descriptions of that term. An ontology used in the data may not be up to date with the latest knowledge evolution of a domain, but if it is archived and accessible, then data users outside an institution can be aware of the intended meaning of a concept by referring to the ontology. Nevertheless, ways and technologies of packaging ontologies and vocabularies with data archives may also be discussed and coordinated, otherwise various bottom-up approaches will lead to additional heterogeneities among data sources. This can be a topic for discussion in the communities of practice discussed above.

The third recommendation is to seek methods for automatic or semi-automatic reworking of data in a Semantic Web context. The Semantic Web not only provides new formats for publishing data online but also leverages the functionalities of querying and inferencing, which streamline the processes in a data life cycle and enables dynamic data integration and reworking. A part of geoscience data such as spreadsheets and relational databases can be converted into Semantic Web formats such as the Resource Description Framework (RDF) (Han et al., 2008). Works have been done towards automatic reworking of data following ontology evolution (Kondylakis and Plexousakis, 2011). Nonetheless, to realize quick updates of online spatial data like digital maps, either triggered by a dynamical change in ontology or not, more efforts are still needed, but Semantic Web technologies lay out a platform for innovative approaches. A recent study (Brodaric, 2012) proposes a structure for knowledge evolution in geologic mapping involving the interaction of abduction, induction, and deduction, and further enriches the structure to characterize and represent inference histories in geologic mapping. This structure is also useful to the quick reworking of digital maps. Such work will maintain consistent semantics within a data life cycle and create an environment for the preservation of ontologies and datasets.

4 CONCLUDING REMARKS

Though we have based our discussion and recommendations in this paper in geoscience, we assert that necessarily knowledge evolves in science, engineering, medicine, humanities, etc., and there are significant implications for the encoding of such knowledge in ontologies. Ontology dynamics, such as type/format variation and version changes, can introduce considerable challenges to data management and re-use, as ontologies are deployed in various parts of a data life cycle. Ontology dynamics may lead to: the need for minimal to substantial reworking of the extant data, semantic mismatch among data providers, semantic mismatch between data providers and data users, as well as error propagation in cross-discipline data discovery and re-use. Though some of these challenges may never be fully resolved, we summarize our recommendations of actions to reduce their negative impacts in a data life cycle:

- (1) develop and use communities of practice and collaborative approaches on ontologies;
- (2) use ontologies in the procedure of data collection and make them accessible to data users; and
- (3) seek methods for automatic or semi-automatic reworking of data in a Semantic Web context.

ACKNOWLEDGMENTS

We are thankful to the reviewers of the paper for their insightful and constructive comments. We also want to thank colleagues in the CGI-IUGS and ESIP for discussing various topics related to geoscience ontologies and vocabularies.

REFERENCES CITED

- Allen, D., Schuurman, N., Deshpande, A., et al., 2008. Data Integration and Standardization in Cross-Border Hydrogeological Studies: A Novel Approach to Hydrostratigraphic Model Development. *Environmental Geology*, 53(7): 1441–1453
- Bates, R. L., Jackson, J. A., 1995. Glossary of Geology, 3rd Edition. American Geological Institute, Alexandria, VA. 788
- Bell, D., Qi, G., Liu, W., 2007. Approaches to Inconsistency Handling in Description-Logic Based Ontologies. In: Meersman, R., Tari, Z., Herrero, P., eds., On the Move to Meaningful Internet Systems 2007: Otm 2007 Workshops, Pt 2, Proceedings. Springer-Verlag, Berlin. 1303–1311
- British Geological Survey, 2012. British Geological Survey Taxonomy Online. <https://www.bgs.ac.uk/taxonomy/home.html>. Accessed on July 08, 2012
- Brodaric, B., 2004. The Design of GSC FieldLog: Ontology-Based Software for Computer Aided Geological Field Mapping. *Computers & Geosciences*, 30(1): 5–20
- Brodaric, B., 2012. Characterizing and Representing Inference Histories in Geologic Mapping. *International Journal of Geographical Information Science*, 26(2): 265–281
- CCOP, CIFEG, 2006. Asian Multilingual Thesaurus of Geosciences. Coordinating Committee for Geoscience Programmes in East and Southeast Asia (CCOP), Bangkok, Thailand and Centre International pour la Formation et les Echanges en Géosciences (CIFEG), Orléans. 563
- Cox, S. J., 2011. OWL Representation of the Geologic Time-scale Implementing Stratigraphic Best Practice. Proceedings of AGU 2011 Fall Meeting, San Francisco. Abstract IN31B-1440
- Cox, S. J., 2012. Vocabularies of Geologic Time Scale. <http://resource.geosciml.org/vocabulary/timescale>. Accessed on September 6, 2012
- Cox, S. J. D., Richard, S. M., 2005. A Formal Model for the Geologic Time Scale and Global Stratotype Section and Point, Compatible with Geospatial Information Transfer Standards. *Geosphere*, 1(3): 119–137
- Flouris, G., D'Aquin, M., Antoniou, G., et al., 2009. Special Issue on Ontology Dynamics. *Journal of Logic and Computation*, 19(5): 717–719
- Flouris, G., Manakanatas, D., Kondylakis, H., et al., 2008. Ontology Change: Classification and Survey. *The Knowledge Engineering Review*, 23(2): 117–152
- Geo-Data Informatics Workshop Committee, 2011. NSF Geo-

- Data Informatics: Exploring the Life Cycle, Citation and Integration of Geo-Data (Workshop Report). http://tw.rpi.edu/media/latest/WorkshopReport_GeoData2011.pdf. Accessed on July 4, 2012
- Gruber, T. R., 1995. Toward Principles for the Design of Ontologies Used for Knowledge Sharing. *International Journal of Human-Computer Studies*, 43(5–6): 907–928
- Guarino, N., 1997. Understanding, Building and Using Ontologies. *International Journal of Human-Computer Studies*, 46(2–3): 293–310
- Han, L., Finin, T., Parr, C., et al., 2008. RDF123: from Spreadsheets to RDF. In: Sheth, A., Staab, S., Dean, M., et al., eds., *The Semantic Web-ISWC 2008*, LNCS vol. 5318. Springer-Verlag, Berlin. 451–466
- International Commission on Stratigraphy, 2012. Global Boundary Stratotype Section and Point (GSSP) of the International Commission of Stratigraphy. <http://www.stratigraphy.org/column.php?id=GSSPs>. Accessed on July 31, 2012
- Janowicz, K., Hitzler, P., 2012. The Digital Earth as Knowledge Engine. *Semantic Web*, 3(3): 213–221
- Kondylakis, H., Plexousakis, D., 2011. Ontology Evolution in Data Integration: Query Rewriting to the Rescue. Jeusfeld, M., Delcambre, L., Ling, T. W., eds., *Conceptual Modeling-ER 2011*, LNCS vol., 6998. Springer-Verlag, Berlin. 393–401
- Laxton, J., Serrano, J. J., Tellez-Arenas, A., 2010. Geological Applications Using Geospatial Standards—An Example from OneGeology-Europe and GeoSciML. *International Journal of Digital Earth*, 3(Suppl.): 31–49
- Ma, X., Asch, K., Laxton, J. L., et al., 2011a. Data Exchange Facilitated. *Nature Geoscience*, 4(12): 814
- Ma, X., Carranza, E. J. M., Wu, C., et al., 2011b. A SKOS-Based Multilingual Thesaurus of Geological Time Scale for Interoperability of Online Geological Maps. *Computers & Geosciences*, 37(10): 1602–1615
- Ma, X., Fox, P., 2013. Recent Progress on Geologic Time Ontologies and Considerations for Future Works. *Earth Science Informatics*, 6(1): 31–46 doi:10.1007/s12145-013-0110-x
- Ma, X., Wu, C., Carranza, E. J. M., et al., 2010. Development of a Controlled Vocabulary for Semantic Interoperability of Mineral Exploration Geodata for Mining Projects. *Computers & Geosciences*, 36(12): 1512–1522
- Mascarelli, A. L., 2009. Quaternary Geologists Win Timescale Vote. *Nature*, 459: 624
- McGuinness, D. L., 2003. Ontologies Come of Age. Fensel, D., Hendler, J., Lieberman, H., et al., eds., *Spinning the Semantic Web: Bringing the World Wide Web to Its Full Potential*. MIT Press, Cambridge. 171–196
- NADM Steering Committee, 2004. NADM Conceptual Model 1.0—A conceptual Model For Geologic Map Information: U.S. Geological Survey Open-File Report 2004-1334, North American Geologic Map Data Model (NADM) Steering Committee, Reston. 58
- Noy, N. F., Musen, M. A., 2004. Ontology Versioning in an Ontology Management Framework. *IEEE Intelligent Systems*, 19(4): 6–13
- Obrst, L., 2003. Ontologies for Semantically Interoperable Systems. Proceedings of the Twelfth International Conference on Information and Knowledge Management, New Orleans. 366–369
- Raskin, R. G., Pan, M. J., 2005. Knowledge Representation in the Semantic Web for Earth and Environmental Terminology (SWEET). *Computers & Geosciences*, 31(9): 1119–1125
- Sen, M., Duffy, T., 2005. GeoSciML: Development of a Generic GeoScience Markup Language. *Computers & Geosciences*, 31(9): 1095–1103
- Soller, D., Berg, T., 2005. The U.S. National Geologic Map Database Project: Overview & Progress. Ostaficzuk, S. R., ed., *The Current Role of Geological Mapping in Geosciences*. Springer, Dordrecht. 245–277
- Spencer, S., 2012. What is DDI? <http://www.ddialliance.org/what>. Accessed on July 4, 2012

APPENDIX

Listings

isc:GSSPBaseQuaternary

rdf:type gts:GlobalStratotypePoint;

rdfs:label "Location of GSSP Base Quaternary";

gts:primaryGuidingCriterion

"Vrica, Italy | 39.0385 N 17.1348 E | base of the marine claystone overlying the sapropelic marker Bed e (Mediterranean Precession Related Sapropel, MPRS 176) | Magnetic 15 kyr after end of Olduvai (C2n) normal polarity chron | Episodes 8/2, p. 116 120, 1985"@en;

gts:status "Ratified 1985"@en.

Listing 1. Definition of GSSP of the Quaternary base for the 2008 International Stratigraphic Chart in Turtle syntax (Cox, 2012).

isc:GSSPBaseQuaternary

rdf:type gts:GlobalStratotypePoint;

rdfs:label "Location of GSSP of Base Quaternary"@en;

gts:primaryGuidingCriterion

"Monte San Nicola, Sicily, Italy | 37.1469 N 14.2035 E | base of marly layer overlying sapropel MPRS 250 | Magnetic Matuyama/Gauss boundary (C2r/C2An) is 1m below GSSP. GSSP level is within Marine Isotope Stage 103. | Episodes 21/2, p. 82 87, 1998"@en;

gts:status "Ratified 1996"@en.

Listing 2. Definition of GSSP of the Quaternary base for the 2009 International Stratigraphic Chart in Turtle syntax (Cox, 2012).