



This is the final version of the following article:

León Araúz, Pilar, and Pedro Javier Magaña Redondo. 2010. EcoLexicon: contextualizing an environmental ontology. In *Proceedings of the Terminology and Knowledge Engineering (TKE) Conference 2010*. Dublin: Dublin City University.

You can find more articles authored by LexiCon Research Group members at <<http://lexicon.ugr.es>>.

EcoLexicon: contextualizing an environmental ontology

Pilar León Araúz, Pedro Javier Magaña Redondo
University of Granada

1 Introduction

EcoLexicon¹ is a Terminological Knowledge Base (TKB) on environment enhanced by both linguistic and knowledge representation techniques. Our TKB is primarily hosted in a relational database (RDB) but at the same time integrated in an ontological model. TKBs can find in ontologies a powerful representational model, as they add the semantic expressiveness lacking in RDBs. This enables potential queries to be richer, since reasoning techniques can be applied to extract implicit information. In turn, the design of ontologies can also benefit from the theoretical background of linguistics, especially from cognitive approaches.

Our TKB is structured around an Environmental Event (EE) which provides the conceptual underpinnings for the location of conceptual sub-hierarchies (Faber et al. 2006).

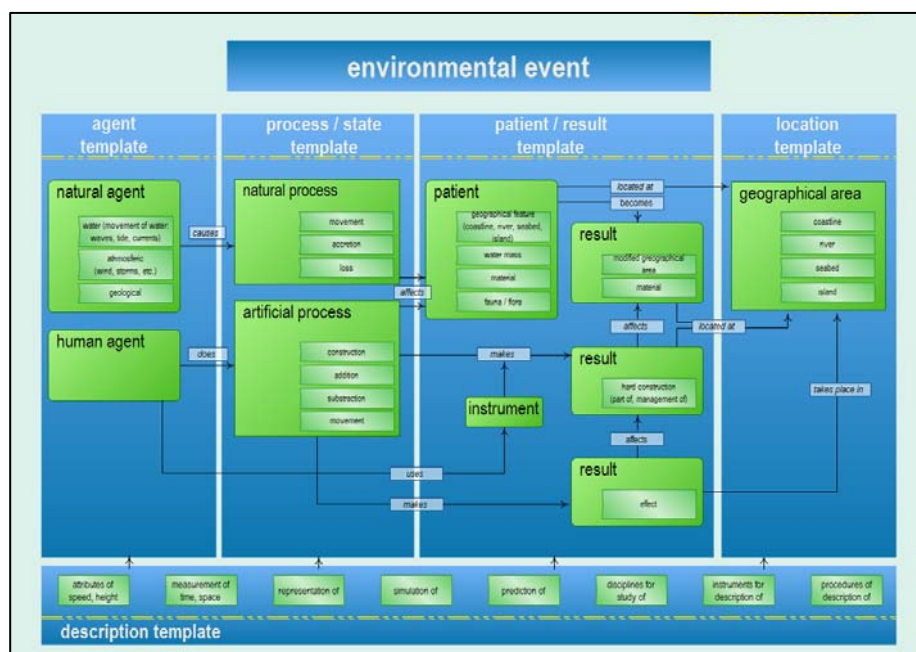


Fig. 1 – Environmental Event

¹ <http://ecolexicon.ugr.es/>

The EE is based on the cognitive linguistics view of frames and semantic roles. . Fillmore and Atkins (1992) define frames as a network of concepts related in such a way that one concept evokes the entire system. According to our corpus-based analysis (Faber et al., 2006), the underlying structure of the entire environmental domain can be encoded in various prototypical frames. Consequently, the upper-level classes in our ontology correspond to the basic semantic roles described in the EE (AGENT-PROCESS-PATIENT-RESULT-LOCATION), all derived from a general knowledge hierarchy (Fig. 2). This structure enables users to gain a better understanding of the complexity of environmental events, since they give a process-oriented general overview of the domain.

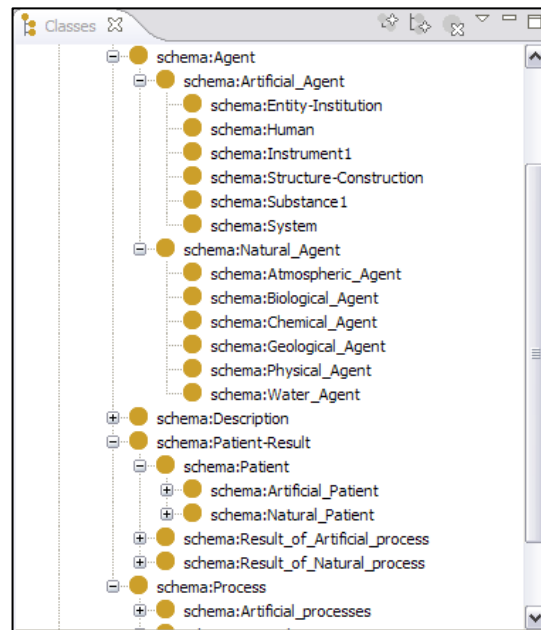


Fig. 2 – Ontological classes

On the other hand, ontologies provide a suitable schema for sharing and reusing semantic resources. According to the Semantic Web initiative, our TKB can benefit from previous works in this field. This could enrich our system with new information, complementing our TKB from a different perspective or even with other contents, such as real-world geographical instances. Nevertheless, information overload not only occurs when interconnecting different systems. Before considering the interoperability of other environmental knowledge-based projects, we must first deal with overinformation in our own TKB.

2 Contextual reconceptualization in EcoLexicon

The final aim of EcoLexicon is to guide the knowledge acquisition process of end users, both for communicative and cognitive purposes. This involves the design of a user-friendly interface where concepts are related in a meaningful way. Based on the EE, conceptual networks in EcoLexicon are structured around a set of different vertical and horizontal

relations, some of which are domain-specific. However, such a wide domain as the ENVIRONMENT has caused an information overload:



Fig. 3 – Information overload

Obviously, users would not acquire any meaningful knowledge if all dimensions of WATER were shown at the same time, as in Figure 3. Overinformation results from a high degree of multidimensionality, which is especially prevalent in what we call *versatile concepts*. Versatile concepts, as WATER, are usually general concepts involved in a myriad of events. For instance, in Figure 3, WATER is linked to the same extent to diverse natural and artificial processes, such as EROSION or DESALINATION. However, WATER rarely activates, if ever, those relations at the same time, as they evoke completely different situations, where WATER is an *agent* in the first one and a *patient* in the second one.

When it comes to hyponymy, the incompatibility among conceptual facets is even more outstanding. Multidimensionality can occur at an intracategorical level, based on the internal structure of concepts. This means that a concept may be classified according to different perspectives but still in the same context, causing the well-known phenomenon of multiple inheritance. For example, different dimensions, like *salinity* or *location* can be hydrological parameters to classify diverse WATER subtypes, such as FRESH WATER, BRACKISH WATER and SALT WATER or SURFACE WATER and GROUNDWATER. These dimensions are compatible enough to share the same conceptual network, because they all describe physical properties of WATER. Moreover, both dimensions are also related, since FRESH, BRACKISH or SALT WATERS can be at the same time either SURFACE or GROUNDWATER.

Nevertheless, hyponymic dimensions show a different nature depending on the external situations where a concept may appear. In that sense, even though WATER subtypes like PRECIPITABLE WATER, DRINKING WATER and NAVIGABLE WATER represent the same dimension *function*, they are not strict coordinate concepts. They only share the same hyperonym, but they will never evoke a common scene. In this line, Barsalou (2005) states that a given concept produces many different situated conceptualizations, each tailored to different instances in different settings. Thus, context can be said to be a dynamic construct that triggers or restricts knowledge.

Our claim is that any specialized domain contains sub-domains in which conceptual dimensions become more or less salient depending on the activation of specific contexts. This means that concepts tend to be intermingled in different situations, and they are not always related to the same concepts or through the same relations. As a result, a more believable representational system should account for re-conceptualization according to the situated nature of concepts.

domains can form their own hierarchical structure. Moreover, they are also dynamic and flexible structures that should evolve over time according to the type and amount of information stored in our TKB. If many other concepts were added to a particular context, new constraints should be developed in accordance with other versatile concepts' special needs. Dynamism would thus help to avoid potential overinformation caused by new data.

On the other hand, overinformation can affect conceptual networks at different stages: (1) when the first concept is versatile, as shown above (2) and when specific concepts are linked to one of the versatile concepts in the first hierarchical level, since the second level will spread the whole first-level network of the versatile concept.

This is why contextual constraints have been applied at all levels, reconceptualising any concept somehow linked to versatile ones, whether they show an information overload or not. In this way, the context-free network of EROSION in figure 6 becomes restricted by the HYDROLOGY context domain shown in figure 7.

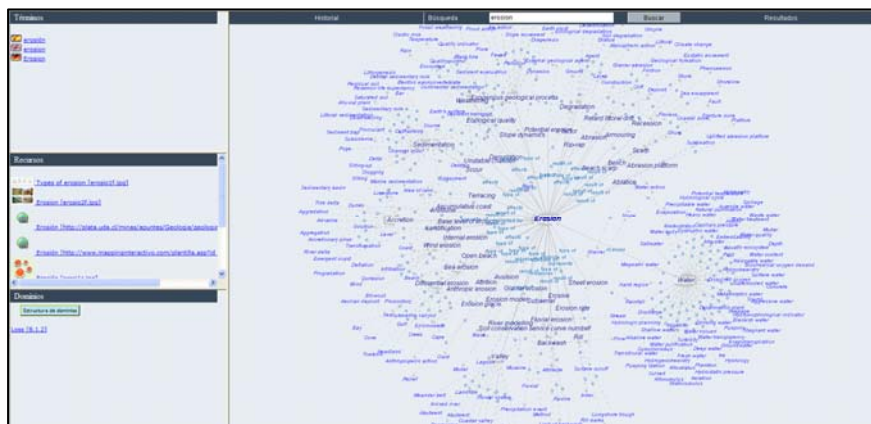


Fig. 6 – EROSION in the context-free network

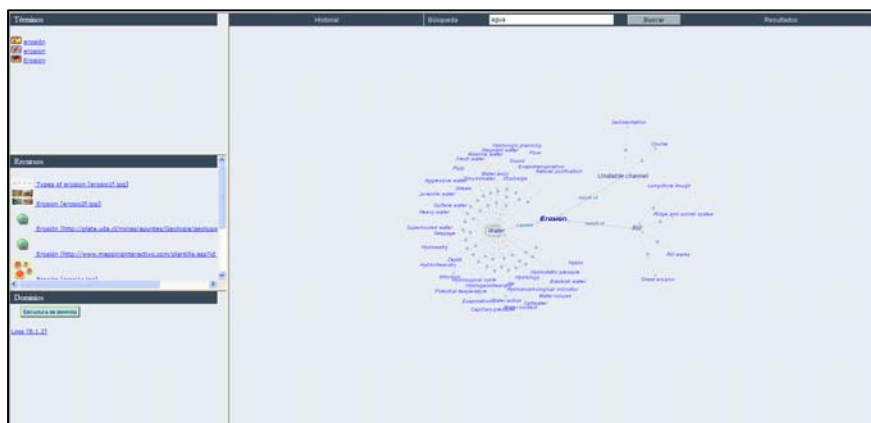


Fig. 7 – EROSION in the HYDROLOGY context domain

3 “Contextualizing” the environment across the web

As mentioned above, our domain knowledge is represented by using a relational database. This widespread modeling let us do a quick deployment of the platform and feed the system from very early stages. Nevertheless, relational modeling has some limitations. One of the biggest ones is its limited capability to represent real-world entities. Ontologies arose as an excellent alternative, but keeping all the development carried out so far was our priority. This is why we emphasize the importance of storing semantic information in the ontology, while leaving the rest in the relational database. In this way, we can continue using the new ontological system, while at the same time feeding the database.

Once information can be accessed by using ontological resources, it is easier to connect it with other environmental resources. Reusability is often based on data merging, but that would lead to a heterogeneous blending of diverse data founded on very different aims. Linked data is an innovative approach facing this problem. It uses Semantic Web technologies to publish structured data and, at the same time, set links different data sources but keeping them as independent resources.

Nevertheless, this is not an easy task, since both representational models have remarkable differences. In contrast to relational databases, ontologies are highly expressive relational structures where concepts are described in very similar terms to those used by humans. Thus, relational models are suited to organize data structure and integrity, whereas ontologies try to specify the meaning of their underlying conceptualization (Barrasa, 2007).

Our ontological classes are fed through the extraction of stored information in the database. This is done by using the D2RQ tool, which provides a usage scenario where relational databases are maintained as non-legacy applications (Bizer and Seaborne, 2004). D2RQ is a declarative language to describe mappings between both systems. Moreover, these mappings can be conditional, which allows for feeding every class just with its corresponding instances.

The next step in our development is to connect this environmental resource with other resources within the same domain. Several techniques have been proposed so far. Former attempts deal with automatic mediation algorithms in order to map and merge between ontology schemas (de Bruijn et al., 2006). However, a remarkable drawback is that the schemas do not always remain public. On the other hand, many systems provide interfaces to interact with their structured data, the well-known APIs (Application Programming Interfaces). This fact has enabled many developers to combine information from different data sources creating new services known as mashups (Zang et al., 2008). Nevertheless, APIs have several disadvantages. Most of the interfaces are proprietary and it is not possible to set links between data objects.

The Linked Data approach (Berners-Lee, 2006) provides an efficient mechanism to publish structured information on the web while object data from different data sources can be linked at the same time. This is why we think this methodology can be applied with success in EcoLexicon and other data sources in order to create an environmental community within the Linked Data framework. EnvO (Morrison, 2009) and SWEET (Raskin, 2003) ontologies are especially interesting to us. SWEET provides a common semantic framework for various Earth science initiatives whereas EnvO aims at developing a common annotation system for any record in the web community that has an environmental component.

This way, we should be able to have statements like the following in the near future:

```
<http://manila.ugr.es/resource/water>  
owl:sameAs  
<http://purl.org/obo/owl/ENVO#ENVO_00002006>
```

This means that water in EcoLexicon (<http://manila.ugr.es/resource/water>) would be related to the same concept (expressed as ENVO 00002006) in EnvO (http://purl.org/obo/owl/ENVO#ENVO_00002006), enriching our conceptualization with any other new data included in these resources. In this way, other resources can equally enhance their systems with our information, which would help to build a real community of shared data.

4 Conclusions

Contextual constraints enrich the system from both a qualitative and quantitative standpoint. On the one hand, they structure knowledge in a similar way to how things relate in the real world, as well as in the human conceptual system. On the other hand, conceptual dimensions are noticeably reduced with a coherent and consistent method based on a cognitive approach. As a result, the situated representation of versatile concepts is a viable solution for managing overinformation and at the same time enhancing knowledge acquisition processes.

We have proven that a legacy system like EcoLexicon can be integrated in the semantic web. Thanks to this achievement, TKBs can also be linked to other resources through new semantic web technologies like linked data. This step is not concluded yet. In the near future we plan to link EcoLexicon to EnvO and Sweet ontologies extensively. However, the success of this approach will largely depend on the proliferation of other shared initiatives.

References

- BARRASA, J. (2007). Modelo para la Definición Automática de Correspondencias Semánticas entre Ontologías y Modelos Relacionales (PhD dissertation, Universidad Politécnica de Madrid).
- BARSALOU, L.W. (2005). Situated conceptualization. In H. Cohen. & C. Lefebvre. Eds. *Handbook of Categorization in Cognitive Science* p. 619-650. St. Louis.
- BERNERS-LEE, T. (2006). Linked Data. W3C Design Issues.
- BIZER, C. & SEABORNE, A. (2004). D2RQ-Treating Non-RDF Databases as Virtual RDF Graphs, *Proceedings of the 3rd International Semantic Web Conference (ISWC2004)*.
- DE BRUIJN, J., EHRIG, M., FEIER, C., MARTÍN-RECUERDA, F., SCHARFFE, F., AND WEITEN, M. (2006). Ontology Mediation, Merging, and Aligning.
- FABER, P., MONTERO MARTÍNEZ, S., CASTRO PRIETO, M.R., SENSO RUIZ, J., PRIETO VELASCO, J.A., LEÓN ARAÚZ, P., MÁRQUEZ LINARES, C., VEGA EXPÓSITO, M. (2006). Process-oriented terminology management in the domain of Coastal Engineering, *Terminology* 12: 2, p.189-213.
- FILLMORE, C.J., ATKINS, B.T.S. (1992). Towards a frame-based lexicon: the semantics of risk and its neighbours. In A. LEHRER & E. F. KITTAY. Eds. *Frames, Fields and Contrasts*, Hillsdale, New Jersey: Lawrence Erlbaum Associates. p. 75-102.
- MORRISON, N. (2009). EnvO - Development of an Environmental Ontology. *Proceedings of Towards eEnvironment. Opportunities of SEIS and SISE: Integrating Environmental Knowledge in Europe*.
- RASKIN, R. & PAN, M. (2003). Semantic Web for Earth and Environmental Terminology (SWEET). *Proceedings of the Workshop on Semantic Web Technologies for Searching and Retrieving Scientific Data*.
- ZANG, N., ROSSON, M. B., AND NASSER, V. (2008). Mashups: who? what? why? In *CHI '08: CHI '08 extended abstracts on Human factors in computing systems*, p. 3171-3176.