

Knowledge Extraction and Representation: the EcoLexicon Methodology

Pilar León Araúz, Arianne Reimerink

Department of Translation and Interpreting, University of Granada

Buenuceso 11, 18002, Granada, Spain

E-mail: pleon@ugr.es, arianne@ugr.es

Abstract

EcoLexicon, a multilingual terminological knowledge base (TKB) on the environment, provides an internally coherent information system which aims at covering a wide range of specialized linguistic and conceptual needs. Knowledge is extracted through corpus analysis. Then it is represented and contextualized in several dynamic and interrelated information modules. This methodology solves two challenges derived from multidimensionality: 1) it offers a qualitative criterion to represent specialized concepts according to recent research on situated cognition (Barsalou, 2009), and 2) it is a quantitative and efficient solution to the problem of information overload.

Keywords: knowledge extraction, knowledge representation, EcoLexicon, multidimensionality, context

1. Introduction

EcoLexicon¹ is a multilingual knowledge base on the environment. So far it has 3,283 concepts and 14,695 terms in Spanish, English and German. Currently, two more languages are being added: Modern Greek and Russian. It is aimed at users such as translators, technical writers, environmental experts, etc., which can access it through a friendly visual interface with different modules devoted to both conceptual, linguistic, and graphical information.

In this paper, we will focus on some of the steps applied to extract and represent conceptual knowledge in EcoLexicon. According to Meyer et al. (1992), terminological knowledge bases (TKBs) should reflect conceptual structures in a similar way to how concepts relate in the human mind. The organization of semantic information in the brain should thus underlie any theoretical assumption concerning the retrieval and acquisition of specialized knowledge concepts as well as the design of specialized knowledge resources (Faber, 2010). In Section 2, we explain how knowledge is extracted through corpus analysis. In Section 3, we show how conceptual knowledge is represented and contextualized in dynamic and interrelated networks.

2. Conceptual Knowledge Extraction

According to corpus-based studies, when a term is studied in its linguistic context, information about its meaning and its use can be extracted (Meyer & Mackintosh, 1996). In EcoLexicon, the corpus consists of specialized (e.g. scientific journal articles, thesis, etc.), semi-specialized texts (textbooks, manuals, etc.) and texts for the general public, all in the multidisciplinary domain of the environment. Each language has a separate corpus and the knowledge is extracted bottom-up from each of the corpora. The underlying ontology is language independent and based on the knowledge extracted from all the corpora. The extraction of conceptual knowledge combines direct term searches and knowledge pattern (KP) analysis. According to many studies on the subject, KPs are considered one of the most reliable methods for knowledge extraction (Barrière, 2004). Normally, the most recurrent knowledge patterns (KPs) for each conceptual relation identified in previous research are used to find related term pairs (Auger & Barrière, 2008). Afterwards, these terms are used for direct term searches to find new KPs and relations. Therefore, the methodology consists of the cyclic repetition of both procedures.

When searching for the term EROSION, conceptual concordances show how different KPs convey different

¹ <http://ecolexicon.ugr.es>

relations with other specialized concepts. The main relations are *caused_by*, *affects*, *has_location* and *has_result*, which highlight the procedural nature of the concept and the important role played by non-hierarchical relations.

In Figure 1, EROSION is related to its diverse kinds of

agents, such as STORM SURGE (1, 7), WAVE ACTION (2, 13), RAIN (3), CONSTRUCTION PROJECTS (6) and HUMAN-INDUCED FACTORS (11). They can be retrieved thanks to all KPs expressing the relation *caused_by*, such as *resultant* (1), *agent for* (2, 3), *due to* (6, 7), and *responsible for* (11).

```

Caused_by
1 Alabama. Significant storm surge and resultant beach erosion were associated with Ivan's landfall. However,
2 nd climate on the Castellon coast, the main agent for erosion is wave action, and this is therefore responsi
3 f a stream. The first factor, rain, is the agent for erosion, but the degree of erosion is governed by oth
4 rts (SW) and semiarid steppe (BS). Wind can also cause erosion and deposition in environments where sediments
5 ety. Reflection of waves from a jetty may also cause erosion of adjacent shorelines. However, erosion furthe
6 ostial zone management. However, in some cases coastal erosion can be due to construction projects that a
7 tude of about 0.3 M m3 per year. Acute erosion Acute erosion due to storm surges (waves and water levels at
8 er. Mangrove removal is also reported to cause coastal erosion and change sedimentation patterns and shoreline
9 [edit] Erosion surface runoff is one of the causes of erosion of the earth's surface. Reduced crop product
10 pes. Local disturbances, for instance by flood-induced erosion, redistribution of sediment or accumulation of
11 ors and human-induced factors, responsible for coastal erosion and highlight the time and space patterns withi
12 cess is typical of a cyclical process of storm-caused erosion in winter, followed by progradation owing
13 can cause excessive wave action that can lead to beach erosion. Trash dumped from boats can be washed up onto
14 that have reached base level develop broad valleys by erosion caused by meandering channels. The stream chain

Affects
15 ing these sensitive creatures. In some cases, coastal erosion can have adverse effects on water quality and h
16 use of dredged material to restore beaches damaged by erosion. EPA works with the U.S. Coast Guard to regul
17 reasonable points, though when push comes to shove and erosion threatens buildings, traditional beach maintena
18 ks and arches found on irregular rocky coastlines; and erosion provides the material which forms deltas and ba
19 near the base of the cliff. This process undercuts and erosion causes the cliffs to retreat landward.

Has_location
21 ed by the position of sand accumulation and beach erosion around littoral barriers. A coastal structure i
22 hes. Kuenen (1950) estimates that beach and cliff erosion along all coasts of the world totals about 0.12g
23 ce and divergence of wave energy over an offshore bar. erosion downdrift of a structure such as a groin, sudd
24 proportional to the longshore transport rate, and erosion takes place downdrift at about the same rate. T

Has_result
25 Excessive loads of silt and other sediments caused by erosion can suffocate bottom-dwelling plants and animal
26 islands or coral reefs. Primary coasts are created by erosion (the wearing away of soil or rock), deposition
27 \par transported. Beach material is also derived from erosion of the coastal formations caused by waves
28 ed to the passage of the ice. Shorelines produced by erosion of glacial till deposits differ markedly from
29 beaches and marshes, are being formed as a result of erosion also transportation of unconsolidated material
30 ion of the seashore and a rise in s.l.w. The results of erosion could lead to further seawater intrusion that c
31 fs are deposited in landslide debris. In this cliffs, erosion of softer material has created bays. The expect
32 s of steep systems, a sea-level rise may cause coastal erosion resulting in profile steepening, and therefore
    
```

Figure 1: Non-hierarchical relations associated with EROSION

```

Is_a
33 vided by the area (A) of the drainage basin (L) Erosion is the natural process of removal of soil by wa
34 in the Netherlands, geomorphological processes such as erosion, transport and sedimentation of sandy materials
35 BURY AND DUXBURY, 1996). coastal processes such as erosion and accretion are site-specific, season-specific
36 these catchments include: stormwater impacts such as erosion, channelisation, sediment deposition and sediment

Type_of
37 eroded by shallow overland flow (sheet, rill and gully erosion) and delivered to the drainage network. Channel
38 m the great local relief, the result of differential erosion by glacier ice. Figure 9-20 includes two sche
39 ing flood events, the dikes are subject to the lateral erosion of the river trying to reoccupy its former coa
40 d enlarges these small channels and generates headward erosion directed towards the aggrading active channel (
41 out five percent of the material on most beaches, wave erosion of rocky coasts is usually slow, even where the
42 of the earth's land surface is dominated by fluvial erosion, lakes that do occur are threatened with either
43 ind climate, topography and surface roughness. wind erosion risk applies only when soils are dry and not c
44 oportional to the steepness of the land surface, water erosion is in proportion to the shear stress exerted by
45 lay to become both wider and deeper over time. Glacial erosion also results in a change in the valley's cross-
46 dominate in periglacial environments: nivation; solifluction; and fluvial erosion and deposition; and fluvial erosion and deposit
47 erosion processes. 215 CHAPTER 13 EQUATIONS: SEDIMENT Erosion caused by rainfall and runoff is computed with
48 givers to simulate cross-shore beach, berm, and dune erosion produced by storm waves and water levels. The l
49 uctures constructed to date have resulted in shoreline erosion in their lee. Furthermore, the key environmen
    
```

Figure 2: Hierarchical relations associated with EROSION

This relation can also be conveyed through compound names such as *flood-induced* (10) or *storm-caused* (12) and any expression containing *cause* as a verb or noun: *one of the causes of* (9), *cause* (4, 5, 8) and *caused by* (14). EROSION is also linked to the patients it *affects*, such as WATER (15), SEDIMENTS (16), and BEACHES (17). However, the affected entities, or patients, are often equivalent to locations (eg. if EROSION *affects* BEACHES it actually *takes place at* the BEACH). The difference lies in the kind of KPs linking the propositions. The *affects* relation is often reflected through the preposition *of* (10) or verbs like *threatens* (18), *damaged by* (17) or *provides* (19), whereas the *has_location* relation is conveyed through prepositions linked to directions (*around*, 21; *along*, 22; *downdrift*, 23) or spatial expressions such as *takes place* (24). In this way, EROSION appears linked to the following locations: LITTORAL BARRIERS (21), COASTS (22) and STRUCTURES (23). *Result* is an essential

dimension in the description of any process, since it also has certain effects, which can be the creation of a new entity (SEDIMENTS, 25; MARSHES, 29; BAYS, 31) or the beginning of another process (SEAWATER INTRUSION, 31; PROFILE STEEPENING, 32).

All these related concepts are quite heterogeneous. They belong to different paradigms in terms of category membership or hierarchical range. For instance, some of the agents of EROSION are natural (WIND, WAVE ACTION) or artificial (JETTY, MANGROVE REMOVAL) and others are general concepts (STORM) or very specific (MEANDERING CHANNEL). This explains why knowledge extraction must still be performed manually, but it also illustrates one of the major problems in knowledge representation: multidimensionality (Rogers, 2004).

This is better exemplified in the concordances in Figure 2, since multidimensionality is most often codified in the *is_a* relation. In the scientific discourse community,

