# GEONLP: A TOOL FOR THE EXTRACTION OF SEMANTIC INFORMATION FROM DEFINITIONS

M. Kokla

School of Rural and Surveying Engineering, National Technical University of Athens, 15780 Zografos Campus, Athens, Greece, mkokla@survey.ntua.gr

**Commission II, WG II/6**

**KEY WORDS:** GIS, Interoperability, Information, Extraction, Integration, Representation, Exchange, Pattern

**ABSTRACT:**

The explication of the semantics of geospatial concepts is a crucial research priority which affects various aspects of information representation, formalization, integration, and exchange. The aim of the present paper is twofold. Firstly, it proposes a methodology for the semantic definition of geospatial concepts. The proposal is based on an analysis of the semantics of geospatial concepts described in information sources such as categorizations, ontologies, data standards, lexical databases, etc. The paper proposes the analysis of semantic information into two types: (a) semantic properties and (b) semantic relations, and provides a list of fundamental semantic properties and relations. Secondly, the paper presents a tool for the extraction and formalization of semantic information from geospatial concept definitions. The tool is used to analyze the definition of each concept and extract the semantic properties and relations and their corresponding values that describe the concept. The output may be used for several tasks, such as concept comparison, ontology development and integration, and semantic information representation.

## 1. INTRODUCTION

The realization of interoperability among various information sources necessitates that the meaning of the exchanged information is properly understood among the interoperating parties. Therefore, as long as the technical problems of information exchange (e.g., protocols, languages, architectures) were resolved, the emphasis was put on the semantic issues. Semantic interoperability became a high priority for the geospatial sciences due to the need to reconcile the differences in the conceptualization and representation of geospatial concepts and preserve their meaning. The explication and preservation of the meaning of geospatial concepts facilitates information comparison, integration, exchange, and reuse.

However, there is no consensus concerning the elements, which specify the meaning of geospatial concepts. Furthermore, existing geospatial information sources such as categorizations, nomenclatures, ontologies, etc. describe geospatial concepts using a variety of elements such as terms, definitions, attributes, relations, instances, etc. Different approaches to information comparison, integration, exchange, and reuse take advantage of one or more of these elements to achieve their goals. Most integration approaches use terms and attributes to integrate concepts from different information sources. The projects KRAFT (Visser et al., 1998) and MOMIS (Beneventano and Bergamaschi, 2004) are representative examples of this approach. Rodriguez and Egenhofer (2003) developed a model to compute semantic similarity between geospatial concepts based on the following components: (a) concept terms, (b) semantic relations among concepts (is-a and part-whole), and (c) distinguishing features (i.e., functions, parts and attributes). However, most of the existing geospatial information sources do not include such a wealth of elements to describe their concepts. Duckham and Worboys (2005) developed an approach to information fusion based on algebra and first-order logic which uses instance-level information to infer semantic

information. Tomai and Kavouras (2005) also use instance information to compare the information content of two thematic maps based on channel theory. The schema integration approach developed by the MIGI (Metadata Integration and Geodata Integrity) project (Hakimpour and Geppert, 2002; Hakimpour and Timpf, 2001) compares definitions of concepts in formal ontologies, i.e., definitions specified by logical axioms.

These approaches produce satisfactory results in case geospatial concepts are properly and thoroughly defined according to the elements used (e.g., concept terms, attributes, relations, instances, etc.). However, the majority of existing geospatial information sources, such as ontologies, categorizations, nomenclatures, thesauri, spatial databases, and spatial data standards, specify geospatial concepts using terms and natural language definitions. Other elements that may contribute to an adequate description of concept semantics are either absent or are superficially defined. Functions and parts are examples of elements that are usually absent, whereas attributes and instances are usually not sufficient and reliable to support an integration or comparison endeavor.

The present paper proposes a methodology to extract semantic information from geospatial concept definitions. Definitions are a means of specifying the meaning of concepts and communicating this meaning to others. However, the emergent need to compare, integrate and reuse existing information sources without loss of semantics entails the extraction and formalization of semantic information immanent in definitions.

## 2. DEFINITIONS

The paper focuses on definitions of geospatial concepts in existing categorizations, data standards, ontologies, etc. Definitions are considered to be important sources of general and domain knowledge (Jensen and Binot, 1987; Klavans et al.,

1993; Swartz, 1997). They are widely used for the organization, description, and communication of information. They describe the meaning of concepts using sufficient information in order to differentiate similar concepts and thus they can be further exploited for tasks involving concept comparison, disambiguation, and integration. Definitions use a sublanguage of natural language (Calzolari, 1984) and a special syntax. In contrast to free text, the special structure and content of definitions facilitate the development of tools for the automatic extraction of semantic information. Moreover, definitions preserve the meaning of information and thus ensure the unambiguous interpretation during information exchange and integration.

The scientific field which deals with natural language generation and understanding, human-computer interaction, and information retrieval and extraction is called Natural Language Processing (NLP) and is based on Artificial Intelligence and Computational Linguistics. NLP deals with the representation of knowledge, either general or domain and the association of knowledge representations with linguistic structures such as vocabulary and grammar (Bateman, 1992). Especially the NLP task, which focuses on the extraction of semantic information, is called "semantic information extraction", and is especially relevant to the tasks of ontology development, comparison, and integration.

The notion of *semantic relation* or *semantic role* or *thematic role* refers to the relation of a constituent to the main verb in a clause. Semantic relations may be either general, or domain-specific. The extraction of semantic information is based on the mapping of linguistic expressions with syntactic relations (e.g., subject-verb-object triples) to semantic relations, using several techniques. The primary methodology for the automatic identification of semantic relations is pattern matching (Khoo and Myaeng, 2002). Patterns are words and phrases in definitions systematically used to express specific semantic information. For example, the phrase "[effect] is the result of [cause]" is a pattern which expresses the CAUSE-EFFECT relation.

## 3. SPECIFICATION OF SEMANTIC INFORMATION

In literature, there is no complete list of semantic information that can be extracted from definitions (Barriere, 1997). This is due to the fact that research has been focused more on the identification of hypernyms or IS-A relations and less on other semantic elements. Furthermore, semantic elements may vary according to the dictionary from which they are extracted, or they may be domain-specific. Therefore, in order to specify the semantic elements which are used for the identification of geospatial concepts it was necessary to analyze geospatial concept definitions from existing information sources in order to identify patterns that are systematically used to express specific semantic elements and formulate the appropriate rules for their extraction. Examples of such information sources are geospatial ontologies, standards, and categorizations, such as CYC Upper Level Ontology, WordNet, CORINE Land Cover, DIGEST, SDTS, etc.

While, as mentioned in the previous section, the term *semantic relations* is used in literature to denote the semantic information extracted from free text or from definitions, in order to be explicit a further classification of semantic information found in definitions is pursued in this paper (Kavouras and Kokla,

2008): (a) *semantic properties* refer to internal characteristics of the concept, i.e., characteristics which are formed independently of other concepts, whereas (b) *semantic relations* describe external characteristics, i.e., characteristics which depend on the interaction with other concepts. For example, semantic properties describe information such as PURPOSE, AGENT, SIZE, SHAPE, whereas semantic relations define the IS-A and PART-OF conceptual relations. Furthermore, other geospatially-oriented semantic elements were also identified. For example, properties such as LOCATION, TIME, DURATION, and COVER, as well as relations, such as SURROUNDNESS, ADJACENCY, OVERLAP, DIRECTION, and PROXIMITY are highly relevant to geospatial concepts. The main semantic properties and relation, both general and geospatial, are shown in Tables 1 and 2 respectively.

| SEMANTIC PROPERTIES |
| --- |
| PURPOSE |
| AGENT |
| PROPERTY-DEFINED LOCATION |
| COVER |
| TIME |
| DURATION |
| FREQUENCY |
| SIZE |
| SHAPE |

Table 1. Main semantic properties of geospatial concepts

| SEMANTIC RELATIONS |
| --- |
| IS-A |
| IS-PART-OF |
| HAS-PART |
| RELATIVE POSITION |
| UPWARD VERTICAL RELATIVE POSITION |
| DOWNWARD VERTICAL RELATIVE POSITION |
| IN FRONT OF HORIZONTAL RELATIVE POSITION |
| BEHIND HORIZONTAL RELATIVE POSITION |
| BESIDE HORIZONTAL RELATIVE POSITION |
| SOURCE - DESTINATION |
| SEPERATION |
| ADJACENCY |
| CONNECTIVITY |
| OVERLAP |
| INTERSECTION |
| CONTAINMENT |
| EXCLUSION |
| SURROUNDNESS |
| EXTENSION |
| PROXIMITY |
| DIRECTION |

Table 2. Main semantic relations of geospatial concepts

Specific patterns are systematically used in definitions to denote the above semantic properties and relations. These patterns guide the formulation of the corresponding rules for the extraction of each semantic element and its value. The PURPOSE semantic property is determined by specific phrases containing the preposition "for" (e.g., for (the) purpose(s) of, for, used for, intended for) followed by a noun phrase, present participle, or infinitival clause; the head of the prepositional phrase indicates the value of the semantic property

(Vanderwende, 1995). In the following definition of the concept "dam", a PURPOSE semantic property is identified with the value "preventing flooding":
"dam: a barrier which forms a reservoir for preventing flooding.

The COVER semantic property is identified by specific phrases including the preposition "of" or by phrases such as "covered by" or "covered with". For example, the following definition includes a COVER semantic property with the value "salt water":
"sea: large body of salt water partially enclosed by land".

Adjectives or adjective phrases expressing size such as "large", "small", "big", "of a large volume", and "tall" indicate a SIZE semantic property. For example, the following definitions include a SIZE semantic property:
"river: *large*, natural stream of water"

The head of the noun phrase, which constitutes the definition, most frequently indicates the genus, i.e., the hypernym or IS-A relation, as in the following definitions:
"hotel: a *building* where travelers can pay for lodging and meals and other services"
"hospital: a medical *institution* where sick or injured people are given medical or surgical care".

| SEMANTIC ELEMENTS | EXAMPLE |
|---|---|
| IS-A | **hotel**: a **building** where travellers can pay for lodging and meals and other services |
| LOCATION | **saltpan**: a shallow basin **in a desert region**<br>**watercourse**: natural body of running water flowing **on or under the earth** |
| COVER | **river**: natural stream of **water**, normally of a large volume<br>**body of water**: the part of the earth's surface covered with **water** |
| SIZE | **snowfield**: a permanent **wide** expanse of snow<br>**river**: **large** natural stream of water |
| TIME | **wadi**: gully or streambed in North Africa and the Middle East that remains dry except **during rainy season** |
| PART-OF | **seacoast**: the shore of a sea or ocean |
| SEPERATION | **coastal lagoons**: stretches of salt or brackish water in coastal areas which are separated from the **sea** by a tongue of land or other similar topography |
| SURROUNDNESS | **lake**: body of water surrounded by **land** |

Table 3. Examples of semantic elements found in definitions

The PART-OF semantic relation is identified by prepositional phrases such as "part of" followed by a noun phrase, as in the following definition where the PART-OF relation takes the value "shore or beach":
"foreshore: that part of the shore or beach which lies between the low water mark and the coastline/shoreline"
The HAS-PART semantic relation is specified by phrases such as "consist of", "comprise of", and "composed of". The

following definition contains a HAS-PART semantic relation with "road or path" as the value:
"way: artifact consisting of a road or path affording passage from one place to another".

Geospatial definitions also convey a lot of spatial relations, such as RELATIVE POSITION, TOPOLOGY, PROXIMITY, DIRECTION, etc. Topological semantic relations are further classified into more detailed types, as follows (Kavouras and Kokla, 2008):

- SEPERATION is expressed by phrases such as "separated from",
- ADJACENCY is expressed by phrases such as "adjacent to", "next to",
- CONNECTIVITY is expressed by phrases such as "connected to",
- OVERLAP is expressed by the verb "overlap", or by phrases such as "partly covered by",
- INTERSECTION is expressed by the verb "cross", phrases such as "intersect with" or by prepositional phrases introduced by the prepositions "through" and "via"
- CONTAINMENT is expressed by phrases such as "contained in" or by prepositional phrases introduced by the prepositions "within" and "inside"
- EXCLUSION is expressed by prepositional phrases introduced by the preposition "outside"
- SURROUNDNESS is expressed by phrases such as "surrounded by", "enclosed by", or by prepositional phrases introduced by the prepositions "around", "among", and "between"
- EXTENSION is expressed by verb phrases including the verbs "extend" and "span" or by prepositional phrases introduced by the prepositions "along" and "across".

Table 3 shows some examples of semantic properties and relations in definitions of geospatial concepts.

## 4. GEONLP: FORMALIZATION AND EXTRACTION OF SEMANTIC INFORMATION

The methodology for the automatic extraction and formalization of semantic information from definitions is implemented by a tool developed by the OntoGEO group called GeoNLP (Kokla, 2005; Mourafetis, 2005; Kavouras and Kokla, 2008). It is based on the approach introduced by Jensen and Binot (1987) and further pursued by Vanderwende (1995) and Barriere (1997). The identification, extraction, and formalization of the semantic elements from definitions are based on the special language and syntax of definitions. More specifically, semantic elements are identified on the basis of pattern matching techniques, which map linguistic expressions and their between syntactic relations (e.g., verb, subject, etc.) to semantic elements.

GeoNLP proceeds in two steps: (a) definition parsing and (b) application of pattern matching rules. The first step performs the syntactic analysis of definitions in order to identify the form, function, and syntactic relations of parts of speech. This step is executed based on a tool called DIMAP-4 (CL Research, 2001) developed for the creation and maintenance of dictionaries. GeoNLP exploits the ability offered by DIMAP-4 to parse dictionary definitions. The output of the first step is given in the form of a parse tree. DIMAP-4 analyzes each

definition into its constituent syntactic parts (e.g., noun phrases, prepositional phrases, verb phrases, subjects of a sentence or clause, prepositions, pronouns, verbs, etc.

At the second step, the parsing result is subject to a set of heuristic rules which search for various syntactic and lexical patterns and extract semantic properties and relations and their values. Each semantic element is automatically extracted with a specifically formulated rule. GeoNLP uses its own programming language and allows the user to formulate new rules or to modify the existing ones. This ability is very important for dealing with the semantics of new sources with domain-specific semantic relations and properties, e.g., domain or application ontologies. For example, the rule for the extraction of the semantic property PURPOSE and its value is the following (Vanderwende, 1995):

*If the verb used (intended, etc.) is post-modified by a prepositional phrase with the preposition "for", then there is a PURPOSE semantic property with the head(s) of that prepositional phrase as the value.*

As mentioned in the previous section, the semantic relation HAS-PART is introduced by phrases such as "consist of", "comprised of", "composed of", and "made of". The rule to extract this semantic relation is formulated as follows:

*If the verb consist (comprise, compose, etc.) is post-modified by a prepositional phrase with the preposition "of", then there is a HAS-PART semantic relation with the head(s) of that prepositional phrase as the value.*

GeoNLP offers two possibilities for processing definitions. The first possibility is to load and process an entire ontology, and provide an xml file that contains the values of the semantic elements of each concept. The second possibility is to analyze single definitions for extracting their semantic elements.

Figures 4 and 5 show the interface of GeoNLP for processing single definitions. The upper left window includes the concept term followed by its definition. The right window shows the output of the parsing step, i.e., the parse tree produced by DIMAP-4. The middle left window shows the rule for identifying a specific semantic property or relation. Figure 4 shows the rule for the extraction of the semantic property PURPOSE and Figure 5 the rule for the extraction of the semantic relation SURROUNDNESS. The lower left window shows the value of the semantic property or relation as extracted by the corresponding rule.

The output of GeoNLP is a set of semantic properties and relations and their corresponding values for each geospatial concept processed. This semantic extraction and formalization process may be further exploited for different tasks, such ontology creation and integration, concept comparison, etc. During ontology creation, semantic elements and values may be used: (a) to explicitly define and document geospatial concepts using an ontology editor, or (b) to formalize geospatial concepts using logical axioms. Furthermore, the semantic elements and their values are suitable in cases where explicitness and objectivity are essential, i.e. concept comparison and ontology integration. The explicit and objective representation and formalization of concept semantics provides the basis to compare similar concepts and to create an integrated ontology. Table 6 shows the comparison of homonymous concepts based on the extraction of their semantic elements and values. It is evident that although the concepts "canal" are given the same terms by Ontology A and Ontology B, their definition differs as to the semantic properties NATURE and PURPOSE. This example reveals the complexity of the semantic description of geospatial concepts. Therefore, the two concepts should not be considered equivalent for a subsequent ontology integration task. The formalization of geospatial concepts based on semantic elements and values facilitates concept comparison and reconciliation in the case of ontology integration.

## 5. CONCLUSIONS

The paper describes a methodology and a tool developed in order to extract semantic information from definitions of geospatial concepts in order to exploit the knowledge immanent in them. Definitions are commonly used to describe the semantics of geospatial concepts; most existing information sources are defined based on concept terms and definitions.

The methodology and the tool may be applied to extract semantic information for tasks such as ontology creation, comparison, and integration which need to be performed with maximum objectivity and explicitness and minimum human intervention. The identification and resolution of semantic heterogeneities should not depend on an expert's subjective decisions, or on insufficiently defined features, such as attributes. Since definitions are wealthy sources of semantic information, GeoNLP maybe used in cases where there are no other semantic elements available or when other elements are superficially defined and therefore not suitable in cases where a semantic approach is desirable.

Although a process dealing with semantics could not be fully automated, further systematization is required. For that reason, future plans include the development of: (a) a value processing approach (e.g., identification of synonyms among values of the same semantic elements), and (b) a comparison process based on semantic elements and values.

### REFERENCES

Barriere, C., 1997. From a children's first dictionary to a lexical knowledge base of conceptual graphs, PhD Thesis, School of Computing Science, Simon Fraser University, British Columbia, Canada.

Bateman, J.A., 1992. The Theoretical Status of Ontologies in Natural Language Processing. In Text Representation and Domain Modeling- Ideas from Linguistics and AI, KIT-Report 97 Preuß, S., and Schmitz, B. (eds.), Technische Universitaet Berlin, pp. 50-99.

Beneventano, D. and Bergamaschi, S., 2004. The MOMIS methodology for integrating heterogeneous data sources. Presented at IFIP World Computer Congress, Toulouse France, August 22-27, http://dbgroup.unimo.it/prototipo/paper/ifip2004.pdf (accessed 3 May 2007).

Calzolari, N., 1984. Detecting patterns in a lexical data base. In: *Proc. 22nd Conference on Association for Computational Linguistics*, Stanford, California, pp. 170-173.

Figure 4. Extraction of the semantic property PURPOSE



Figure 5. Extraction of the semantic relation SURROUNDNESS

| | | SEMANTIC ELEMENTS | | | | | | |
| | *GEOSPATIAL CONCEPTS* | IS-A | COVER | PURPOSE | NATURE | SIZE | FLOW | SURROUNDNESS |
|---|---|---|---|---|---|---|---|---|
| **Ontology A** | **stream** | body | fresh water | | natural | | flowing | |
| | **lake** | body | water | | | | | land |
| | **sea** | body | salt water | | | as large as or larger than a lake | | |
| | **canal** | channel | water | irrigation | artificial | | | |
| **Ontology B** | **stream** | body | running water | | natural | | flowing | |
| | **lake** | body | water | | | | | land |
| | **sea** | body | salt water | | | large body | | land |
| | **canal** | channel | water | transportation | natural or artificial | | | |

Table 6. Comparison of homonymous concepts based on their semantic elements and values

Calzolari, N., 1984. Detecting patterns in a lexical data base. In: Proc. 22nd Conference on Association for Computational Linguistics, Stanford, California, pp. 170-173.

CL Research, 2001, DIMAP-4, Dictionary Maintenance Programs, http://www.clres.com (accessed 2 May 2007).

Duckham, M. and Worboys, M.F, 2005. An algebraic approach to automated information fusion. International Journal of Geographic Information Science, 19(5), pp. 537-557.

Hakimpour, F. and Geppert, A., 2002. Global schema generation using formal ontologies. In: Proc. 21st International Conference on Conceptual Modeling ER2002, Springer-Verlag, LNCS Vol. 2503, pp. 307-321.

Hakimpour, F. and Timpf, S., 2002. A step towards geodata integration using formal ontologies. In: Proc. 5th AGILE Conference on Geographic Information Science, Palma, Balearic Islands, Spain.

Jensen, K. and Binot, J.L., 1987. Disambiguating prepositional phrase attachments by using on-line dictionary definitions. Computational Linguistics, 13(3-4), pp. 251-260.

Kavouras, M. and Kokla, M., 2008. Theories of Geographic Concepts: Ontological Approaches to Semantic Integration. CRC Press, Boca Raton, FL, USA.

Khoo, C. and Myaeng, S.H., 2002, Identifying semantic relations in text for information retrieval and information extraction. In The Semantics of Relationships: An Interdisciplinary Perspective, Green, R., Bean, C.A. and Myaeng, S.H. (Eds.), Kluwer Academic Publishers, Netherlands, pp. 161-180.

Klavans, J., Chodorow, M., and Wacholder, N., 1993. Building a knowledge base from parsed definitions. In: Natural Language Processing: The PLNLP Approach , Jensen, K., Heidorn, G., and Richardson, S. (Eds.), Kluwer Academic Publishers, USA, pp. 119-133.

Kokla, M., 2005. Semantic Interoperability in Geographic Information Science, PhD Thesis, National Technical University of Athens, Athens, Greece (in greek).

Mourafetis, G., 2005. Automated extraction and comparison of geographic information from definitions, MSc Thesis, Geoiformatics Postgraduate Course, National Technical University of Athens, Greece (in greek).

Rodríguez, A. and Egenhofer, M., 2003. Determining semantic similarity among entity classes from different ontologies. IEEE Transactions on Knowledge and Data Engineering 15(2), pp. 442-456.

Swartz, N., 1997. Definitions, dictionaries, and meanings, http://www.sfu.ca/philosophy/swartz/definitions.htm (accessed 1 May 2007)

Tomai, E. and Kavouras, M., 2005. Mappings between maps - Association of different thematic contents using situation theory. In Proc. International Cartographic Conference, A Coruna, Spain, July 6-16.

Vanderwende, L., 1995. The analysis of noun sequences using semantic information extracted from on-line dictionaries, Ph.D. thesis, Faculty of the Graduate School of Arts and Sciences, Georgetown University, Washington, D.C.

Visser, P.R.S, Jones, D., Bench-Capon, T.J.M., and Shave, M.J.R., 1998. Assessing heterogeneity by classifying ontology mismatches. In: Formal Ontology in Information Systems, Guarino, N. (Ed.), Amsterdam: IOS Press, pp. 148-162.