

# Ontological resources and question answering

**Roberto Basili (\*), Dorte H. Hansen (\*\*), Patrizia Paggio (\*\*),  
Maria Teresa Pazienza (\*), Fabio Massimo Zanzotto (\*)**

(\*) Dip. di Informatica Sistemi e Produzione  
University of Rome "Tor Vergata"  
{**basili, pazienza, zanzotto**}  
@**info.uniroma2.it**

(\*\*) Centre for Language Technology  
University of Copenhagen  
{**patrizia, dorte**}@**cst.dk**

## Abstract

This paper discusses the possibility of building an ontology-based question answering system in the context of the Semantic Web presenting a proof-of-concept system. The system is under development in the MOSES European Project.

## Introduction

Question Answering (QA) systems (as QA track of the Text Retrieval Conference (TREC-QA) competitions (Voorhees 2001)), are able both to understand questions in natural language and to produce answers in the form of selected paragraphs extracted from very large collections of text. Generally, they are open-domain systems, and do not rely on specialised conceptual knowledge as they use a mixture of statistical techniques and shallow linguistic analysis. Ontological Question Answering systems, e.g. (Woods et al. 1972, Zajac 2000) propose to attack the problem by means of an internal unambiguous knowledge representation. As any knowledge intensive application, ontological QA systems have as intrinsic limitation related to the small scale of the underlying syntactic-semantic models of natural language.

While limitations are well-known, we are still questioning if any improvement has occurred since the development of the first ontological QA system LUNAR. Several important facts have emerged that could influence related research approaches:

- ♦ a growing availability of lexical knowledge bases that model and structure words: WordNet (Miller 1995) and EuroWordNet (Vossen 1998) among others; some open-domain QA systems have proven the usefulness of these resources, e.g. WordNet in the system described in (Harabagiu et al. 2001).
- ♦ the vision of a Web populated by "ontologically" tagged documents which the semantic Web initiative has promoted; in case this vision becomes a reality, it will require a world-wide collaborative work for building interrelated "conceptualisations" of domain specific knowledge
- ♦ the trend in building shallow, modular, and robust natural language processing systems (Abney 1996, Hobbs et al. 1996, Ait-Moktar&Chanod 1997, Basili&Zanzotto 2002) which is making them appealing in the context of ontological QA systems, both for text interpretation (Andreasen et al. 2002) and for database access (Popescu et al. 2003).

Given this background, we are investigating a new approach to ontology-based QA in which users ask questions in natural language to knowledge bases of facts extracted from a federation of Web sites and organised in topic map repositories (Garshol 2003). Our approach is being investigated in the context of EU project MOSES<sup>1</sup>, with the explicit objective of developing an ontology-based methodology to search, create, maintain and adapt semantically structured Web contents according to the vision of the Semantic Web. MOSES is taking advantage of expertise coming from several fields: software agent technology, NLP, graph theory

---

<sup>1</sup> MOSES is a cooperative project under the 5th Framework Programme. The project partners are FINSA Consulting, MONDECA, Centre for Language Technology, University of Copenhagen, University of Roma Tre, University of Roma Tor Vergata and ParaBotS.

and text mining. The test-bed chosen in the project is related to the development of an ontology-based knowledge management system and an ontology-based search engine that will both accept questions and produce answers in natural language for the Web sites of two European universities. The challenges of the project are:

- ◆ building an ontological QA system;
- ◆ developing a multilingual environment which implies the ability to treat several languages, and, importantly, several conceptualisations.

In this paper, after briefly describing how the project is trying to comply with the semantic Web vision, we will focus on question processing, and in particular on the way in which NLP techniques and ontological knowledge interact in order to support questions to specific sites or to site federations.

### An ontology-based approach to question answering

In our ontological QA system, both questions and domain knowledge are represented by the same ontological language. It is foreseen to develop the QA system in two steps. First a prototypical implementation is planned to answer questions related to the current “state-of-affairs” of the site to which the question is posed. In a second step, given a “federation” of sites within the same domain, we will investigate whether and how an ontological approach could support QA across the sites. Answering a question can then be seen as a collaborative task between ontological nodes belonging to the same QA system. Since each node has its own version of the domain ontology, the task of passing a question from node to node may be reduced to a mapping task between (similar) conceptual representations. To make such an approach feasible, a number of difficult problems must still be solved. In this paper, we will provide details on how:

- ◆ to build on existing ontologies and interface between them and language resources;
- ◆ to interpret questions wrt the ontological language;
- ◆ to model the mapping task for federated questions.

### Building on off-the-shelf semantic Web ontologies

One of the results of the Semantic Web initiative will be the production of many interrelated domain-specific ontologies that provide the formal language for describing the content of Web documents. In spite of the freedom allowed in the production of new conceptualisations, it is reasonable to expect that a first knowledge

representation jungle will leave room to a more orderly place where only the more appreciated conceptualisations have survived. This is a prerequisite for achieving interoperability among software agents. In view of this, and since publicly available non-toy ontology examples are already available, the effort of adapting an existing ontology to a specific application is both useful and possible. This experiment is being conducted in MOSES to treat the university domain.

Ontologies for the Semantic Web are written in formal languages (OWL, DAML+OIL, SHOE) that are generalisations/restrictions of Description Logics (Baader et al. 2003). TBox assertions describe concepts and relations. A typical entry for a concept is:

<i>ID</i>	Course
<i>Label</i>	Course
<i>Subclassof</i>	Work

Table 1 A concept

where *ID* is the concept unique identifier, *label* is the readable name of the concept, *subclassof* indicates the relation to another class. As the label has the only purpose of highlighting the concept to human readers, alternative linguistic expressions are not represented. On the contrary, this piece of information is recorded in a lexical data base like WordNet. The problem is even more obvious when considering relationships.

<i>ID</i>	teacherOf
<i>Label</i>	Teaches
<i>Domain</i>	#Faculty
<i>Range</i>	#Course

Table 2 A relationship

In Table 2, *domain* and *range* contain the two concepts related to the described binary relation. The label *teacherOf* does not mention alternative linguistic expressions like: #Faculty gives #Course or #Faculty delivers #Course, etc.

For the ontology producers, only one concept or relation name is sufficient. Synonymy is not a relevant phenomenon in ontological representations. In fact, it is considered a possible generator of unnecessary concept name clashes, i.e. concept name ambiguity. Conceptualisations (as in tables 1,2) are inherently weak whenever used to define linguistic models for NLP applications. Interpreting questions like:

- (1) Who gives/teaches the database class/course this year?

with respect to a university domain ontology means in fact mapping all the questions onto the concepts and relations in Table 2. There is a gap to be filled between linguistic and ontological ways of expressing the domain knowledge.

### *Linguistic interfaces to ontologies*

In developing an ontological QA system, the main problem is to build what we call the “linguistic interface” to the ontology which consists of all the linguistic expressions used to convey concepts and relationships. To make this attempt viable, we are currently studying methods to automatically relate lexical knowledge bases like WordNet (Miller 1995) to domain ontologies (Basili et al 2003a) and to induce syntactic-semantic patterns for relationships (Basili et al 2003b).

The linguistic interface constitutes the basis on which to build the semantic model of the natural language processing sub-system. One way of conceiving such a model is in terms of syntactic-semantic mapping rules that apply to alternative expressions of the same conceptual knowledge. The amount of syntactic analysis such rules foresee will vary according to the approach chosen.

### *Classifying questions*

To facilitate recognition of what are the relevant expressions to be encoded in the linguistic interface, we have introduced a classification of the possible questions that the system is expected to support. A classification often quoted is that in Lauer, Peacock and Graesser (1992), which mainly builds on speech act theory. Another influential, more syntactically-oriented approach is that in Moldovan et al. (1999) where to each syntactic category correspond one or several possible answer types, or focuses (a person, a date, a name, etc.).

Several dimensions have been identified as relevant for MOSES

1. the number of sites and pages in which the answer is to be found. Thus, a first distinction is done between site-specific and federated questions. In the first case, analysis involves only one language and one knowledge domain. In the second, the interpretation of a question produced by a local linguistic analyser is matched against the knowledge domain of other sites;
2. sub-domain coverage (e.g. people, courses, research).
3. format of the answer: which in MOSES is not only a text paragraph as in standard QA, but could also be composed of one or more in-

stances of semantic concepts (professors, courses) or relations (courses being taught by specific professors), whole Web pages, tables, etc. due to the heterogeneity of information sources

These dimensions have been explored in “question cards” defined by the project’s user groups<sup>2</sup>.

FORM 1	
Input	Hvem underviser i filmhistorie ( <i>Who teaches film history</i> )
Syntactic type	Who (Hvem)
Syntactic subtype	V ≠ copula
CONTENT	
Focus constraint	Teacher
Concepts	Faculty Course.Name: <i>history of film</i>
Relations	TeacherOf(Faculty, Course)
Answer count	List

Table 3: Example of question classification

From the point of view of the linguistic analysis, however, syntactic category and content are the central dimensions of sentence classification. Syntactic categories are e.g. *yes/no question*, *what-question*, *who-question*, etc. Subtypes relate to the position inside the question where the focus is expressed, e.g. depending on whether the wh-pronoun is a determiner, or the main verb is a copula. The content consists of concepts and relations from the ontology, the focus constraint<sup>3</sup> (the ontological type being questioned), and a count feature indicating the number of instances to be retrieved. Table 3 shows an example of linguistic classification. For each sentence type, several paraphrases are described.

### **Ontology Mapping in a Multilingual Environment: challenges**

The conceptualisation of the university world as it appears in the DAML+OIL ontology library is an interesting representation for the application scenarios targeted in MOSES (i.e. *People/Course/Research*). Described classes and relations cover in fact, at least at a high level, most of the relevant concepts of the analysed scenarios. Such an ontology has been adapted to develop conceptualisations for each of the two national

<sup>2</sup> The University of Roma III and the Faculty of Humanities at the University of Copenhagen.

<sup>3</sup> In the sense of Rooth (1992).

university sub-systems (i.e. Italian and Danish) while providing additional information required for answering the input questions. This is temporal information or other kind of information at a border line with the domain, (e.g. concepts related to the job market). A first important matter we have dealt with is the language. Whereas concept and relation labels in the Italian ontology are expressed either in English (for concepts directly taken from the original source) or in Italian, in the Danish counterpart all labels are in Danish. This means that a mapping algorithm making use of string similarity measures applied to concept labels will have to work with translation, either directly between the two languages involved, or via a pivot language like English. The goal would be to establish correspondences such as 'Lektor' ↔ ('AssociateProfessor') ↔ 'ProfessoreAssociato'.

Another problem is related to structural differences: not all the nodes in an ontology are represented also in the other and vice-versa, moreover nodes that are somehow equivalent, may have different structural placements. This is the case for the 'Lektor'/'ProfessoreAssociato' pair just mentioned: in the Danish system, 'Lektor' is not a subclass of 'Professor', although "associate professor" is considered a correct translation.

## Question analysis

Question analysis is carried out in the MOSES linguistic module associated with each system node. To adhere to the semantic Web approach, MOSES poses no specific constraints on how the conceptual representation should be produced, nor on the format of the output of each linguistic module. The agent that passes this output to the content matcher (an ontology-based search engine) maps the linguistic representation onto a common MOSES interchange formalism (still in an early development phase). Two independent modules have been developed for Danish and Italian language analysis. They have a similar architecture (both use preprocessing, i.e. POS-tagging and lemmatising, prior to syntactic and semantic analyses), but specific parsers. Whereas the Danish parser, an adapted version of PET (Callmeier 2000) produces typed feature structures (Copestake 2002), the Italian one outputs quasi-logical forms. Both representation types have proven adequate to express the desired conceptual content. As an example, the Italian analysis module is described below.

### Analysis of Italian questions

Analysis of Italian questions is carried out by using two different linguistic interpretation levels. The syntactic interpretation is built by a general purpose robust syntactic analyser, i.e. Chaos (Basili&Zanzotto 2002).

This will produce a Question Quasi-Logical Form (Q-QLF) of an input question based on the extended dependency graph formalism (XDG) introduced in (Basili&Zanzotto 2002). In this formalism, the syntactic model of the sentence is represented via a planar graph where nodes represent constituents and arcs the relationships between them. Constituents produced are chunks, i.e. kernels of verb phrases (VPK), noun phrases (NPK), prepositional phrases (PPK) and adjectival phrases (ADJK). Relations among the constituents represent their grammatical functions: logical subjects (lsubj), logical objects (lobj), and prepositional modifiers. For example, the Q-QLF of the question

- (2) Chi insegna il corso di Database?  
(Who teaches the database course?)

is shown in Figure 1.



Figure 1 A Q-QLF within the XDG formalism

Then a robust semantic analyser, namely the Discourse Interpreter from LaSIE (Humphreys et al. 1996) is applied. An internal world model has been used to represent the way in which the relevant concepts (i.e. objects) and relationships (i.e. events) are associated with linguistic forms (see Figure 2). Under the object node, concepts from the domain concept hierarchy are mapped onto synsets (sets of synonyms) in the linguistic hierarchy EWN (i.e. the EuroWordNet.base concepts). This is to guarantee that linguistic reasoning analysis is made using general linguistic knowledge.

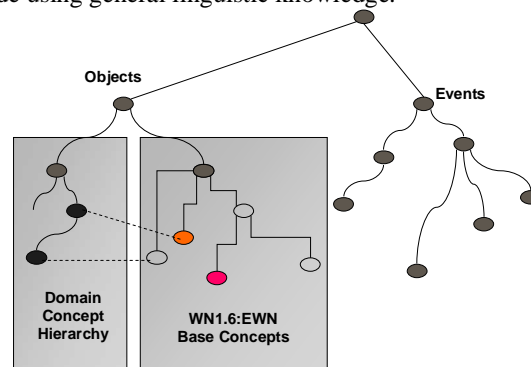


Figure 2 The world model taxonomy

```

TEACH_EVENT ==> teach_course.
teach_course ==> tenere v insegnare v fare.

props(teach_course(E),[
(consequence(E,
[relation(E,teacherOf),r_arg1(E,X),r_arg2(E,Z)] ):-
nodeprop(E,lsubj(E,X)),
X <- ewn4123(_), /* human_1 */
nodeprop(E,lobj(E,Z)),
Z <- ewn567704(_) /* education_1 */
)
]).

```

Figure 3 Example of syntactic-semantic interpretation rule

The association of objects and events with linguistic forms is used in matching rules as shown in Figure 3. The rule expresses the fact that, if any word like *tenere*, *insegnare* or *fare* is encountered in relation with a *human\_1* (represented by the base concept *ewn4123*) and the word *education\_1* (*ewn567704*), the relation *teacherOf* can be induced.

The analysis resulting for sentence (2) is then:

```

focus(e2),
relation(e1,teacherOf),
r_arg1(e1,person_dch(e2)),
r_arg2(e1,course_dch(e3)),
relation(e4,hasSubject),
r_arg1(e4,course_dch(e3)),
r_arg2(e4,topic_dch("Database")).

```

This means that the user is interested in a person, the entity *e2* of the class *person\_dch*, that is in a relation *teacherOf* with the entity *e4* (instance of the class *course\_dch*), that is in turn related by *hasSubject* with the topic (i.e. *topic\_dch*) "Database". This result can be passed on to the content matcher.

### Treating federated questions

Now we want to extend this approach to question analysis in order to manage federated questions. A possible solution would be sending the natural language question to several nodes and let each node interpret it against its own domain knowledge. This is unfeasible in a multilingual environment. The solution we are investigating is based on the notion of ontology mapping. Let us consider the case of a student questioning not only the Danish but also the Italian site (by selecting specific modalities for entering questions):

(3) Hvem er lektor i fransk?

(Who is associate professor of French?)

As the question is in Danish, it has to be analysed by the Danish analysis component, which will produce a semantic interpretation roughly corresponding to the following term:

$\text{all}(x) (\text{lektor}(x) \ \& \ \text{CourseOffer}(x,y) \ \& \ \text{Course}(y) \ \& \ \text{Name}(y, \text{French}))^4$

Since all concepts and relations come from the Danish ontology, it is not a problem to query the Danish knowledge base for all relevant examples. In order to query the Italian knowledge base, however, equivalent concepts and relations must be substituted for those in the "Danish" interpretation. The corresponding Italian representation is:

$\text{all}(x) (\text{ProfessoreAssociato}(x) \ \& \ \text{TeacherOf}(x,y) \ \& \ \text{Course}(y) \ \& \ \text{Subject}(y, \text{French}))$

The first problem is establishing a correspondence between 'lektor' and 'ProfessoreAssociato', which as shown in the ontology fragments below are not structurally equivalent.

As suggested in (Pazienza&Vindigni 2003, Medche&Staab 2001), equivalence relations must be established by considering *is-a* structures and lexical concept labels together. In the example under discussion, an initial equivalence can be posited between the top nodes of the two ontology fragments, since they both refer explicitly to the original DAML+OIL ontology via a *sameAs* relation. However, none of the concept labels under 'Faculty' in the Italian ontology are acceptable translations of 'Lektor', nor do any of the nodes refer to common nodes in a common reference ontology. Thus, the matching algorithm must search further down for an equivalent concept by considering possible translations of concept labels and testing the relations that equivalence candidates participate in. Thus, distance from a common starting node, lexical equivalence and occurrence in similar relations are all constraints to be considered.

<sup>4</sup> All concepts and relations will in fact be expressed in Danish. Here, to facilitate non-Danish readers, we are using English equivalents with the exception of the concept 'Lektor' under discussion.

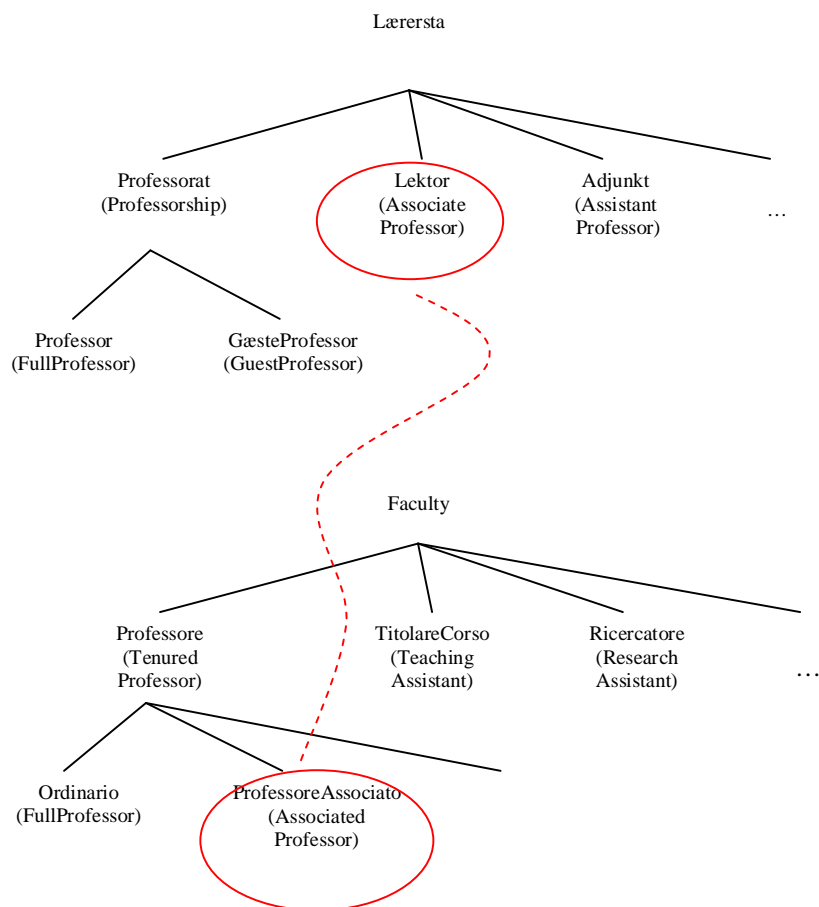


Figure 4: The “Faculty” Danish and Italian sub-ontologies

The same problem of finding a correct mapping appears for the relations. In this case, we must be able to discover that `CourseOffer` and `TeacherOf` represent the same relation. For instance we can rely on the fact that they have both two roles, and the concepts filling these roles, `Faculty` and `Course` (or rather the Danish and Italian equivalent concepts) correspond. Discovering similarities between relations, however, may be a much more complex task than shown in this example. In general, it presupposes the ability to map between concepts.

## Conclusion

Our focus in this paper has been, in the context of ontology-based QA, to discuss how to interface between ontology and linguistic resources on the one hand, and ontology and natural language questions on the other while remaining within a unique framework. An interesting issue in a multilingual environment is how to support questions to federation of sites organised around local ontologies. We have begun to address this issue in

terms of ontology mapping. Specific algorithms for machine learning and information extraction have also been identified and are under development.

## References

Steven Abney (1996) *Part-of-speech tagging and partial parsing*. In G.Bloothoof K.Church, S.Young, editor, *Corpus-based methods in language and speech*. Kluwer academic publishers, Dordrecht.

Salah Ait-Mokhtar and Jean-Pierre Chanod. (1997) *Incremental Finite-state parsing*. In *Proceedings of ANLP97*, Washington.

Andreasen, Troels, Per Anker Jensen, Jørgen F. Nilsson, Patrizia Paggio, Bolette Sandford Pedersen and Hanne Erdman Thomsen (2002) *Ontological Extraction of Content for Text Querying*, in *Natural Language Processing and Information Systems*, Revised Papers of NLDB 2002. Springer-Verlag, pp. 123–136.

- Baader, F., D. Calvanese, D. McGuinness, D. Nardi, P.F. Patel-Schneider, eds. (2003) *The Description Logics Handbook: Theory, Implementation, and Applications*, Cambridge University Press
- Basili, Roberto, Michele Vindigni, Fabio Massimo Zanzotto (2003a) *Integrating ontological and linguistic knowledge for Conceptual Information Extraction*, Web Intelligence Conference, Halifax, Canada, September 2003
- Basili, Roberto, Maria Teresa Pazienza, and Fabio Massimo Zanzotto (2003b) *Exploiting the feature vector model for learning linguistic representations of relational concepts* Workshop on Adaptive Text Extraction and Mining (ATEM 2003) held in conjunction with European Conference on Machine Learning (ECML 2003) Cavtat (Croatia), September 2003
- Basili, Roberto and Fabio Massimo Zanzotto (2002) *Parsing Engineering and Empirical Robustness* Journal of Natural Language Engineering 8/2-3 June 2002
- Burger, John *et al* (2002) *Issues, tasks and program structures to roadmap research in question & answering (Q&A)*. NIST DUC Vision and Roadmap Documents, <http://www-nlpir.nist.gov/projects/duc/roadmapping.html>.
- Callmeier, Ulrich (2000) PET – a platform for experimentation with efficient HPSG processing techniques. In Flickinger, D., Oepen, S., Tsujii, J. and Uszkoreit, H. (eds.) *Natural Language Engineering. Special Issue on Efficient Processing with HPSG*. Vol. 6, Part 1, March 2000, 99–107.
- Copestake, Ann (2002) *Implementing Typed Feature Structure Grammars*. CSLI Publications. Stanford University.
- Garshol, Lars Marius (2003) Living with Topic Maps and RDF. Technical report. <http://www.ontopia.net/topicmaps/materials/tmrdf.html>.
- Harabagiu, Sanda, Dan Moldovan, Marius Păcă, Rada Mihalcea, Mihai Surdeanu, Răzvan Bunescu, Roxana Girju, Vasile Rus, and Paul Morrescu (2001) *The role of lexico-semantic feedback in open-domain textual question-answering*. In Proceedings of the Association for Computational Linguistics, July 2001.
- Hobbs, Jerry R., Douglas E. Appelt, John Bear, David Israel, Megumi Kameyama, Mark Stickel, and Mabry Tyson (1996). *FASTUS: A cascaded finite-state transducer for extracting information from natural-language text*. In Finite State Devices for Natural Language Processing. MIT Press, Cambridge, MA.
- Humphreys, K., R. Gaizauskas, S. Azzam, C. Huyck, B. Mitchell, H. Cunningham, and Y. Wilks (1998) *University of sheffield: Description of the LASIE-II system as used for MUC-7*. In Proceedings of the Seventh Message Understanding Conferences (MUC-7). Morgan Kaufman, 1998.
- Meadche, Alexander and Steffen Staab (2001) *Comparing Ontologies-Similarity Measures and Comparison Study*, Internal Report No. 408, Institute AIFB, University of Karlsruhe, Germany, 2001
- Miller, George A. (1995) WordNet: A lexical database for English. Communications of the ACM, 38(11):39--41, 1995.
- Pazienza, Maria Teresa and Michele Vindigni (2003) *Agent-based Ontological Mediation in IE systems* in M.T. Pazienza ed. Information Extraction in the Web Era, LNAI 2700, Springer Berlin 2003
- Rooth, M. (1992) A Theory of Focus Interpretation. In *Natural Language Semantics*, Vol. 1, No. 1, pp. 75-116.
- Voorhees, Ellen M. (2001) The TREC question answering track. *Natural Language Engineering* 7(4), pp. 361–378.
- Vossen, Piek (1998) *EuroWordNet: A Multilingual Database with Lexical Semantic Networks* Kluwer Academic Publishers, Dordrecht, October 1998
- Woods, W., R. Kaplan, and B. Nash-Weber (1972) *The Lunar Sciences Natural Language Information System: Final Report*. Technical Report, Bolt Beranek and Newman, Number 2378, June 1972.
- Zajac, Remi (2001) *Towards Ontological Question Answering*, ACL-2001 Workshop on Open-Domain Question Answering, Toulouse, France, 2001