

Multilingual Acquisition of Large Scale Knowledge Resources

Adquisición multilingüe de bases de conocimiento de gran escala

Montse Cuadros

Vicomtech-IK4

Mikeltegi Paselekua, 57

2009 Donostia-San Sebastián

mcuadros@vicomtech.org

Resumen: Tesis doctoral en Informática realizada por Montse Cuadros y dirigida por Dr. Lluís Padró y Dr. German Rigau. La defensa de la tesis fue en la facultad de Informática de la Universitat Politècnica de Catalunya el día 22 de noviembre de 2011. El tribunal estuvo formado por Dr. Horacio Rodríguez (Universitat Politècnica de Catalunya), Prof. Dr. Piek Vossen (Vrije Universiteit Amsterdam), Dra. Arantza Díaz de Ilarraza (Universidad del País Vasco), Dra. Irene Castellón (Universitat de Barcelona) y Dr. Roberto Navigli (Sapienza University of Rome), que le otorgaron la nota de Sobresaliente Cum Laude.

Palabras clave: Adquisición de conocimiento, adquisición de léxico, evaluación de recursos, WordNet, desambiguación de acepciones

Abstract: Ph. D. thesis in Computer Science written by Montse Cuadros under the supervision of Dr. Lluís Padró and Dr. German Rigau. The thesis defense was done on 22th November 2011 at the Computer Science Faculty of the Universitat Politècnica de Catalunya. The Doctoral Examination Committee was composed by Dr. Horacio Rodríguez (Universitat Politècnica de Catalunya), Prof. Dr. Piek Vossen (Vrije Universiteit Amsterdam), Dra. Arantza Diaz de Ilarraza (Universidad del País Vasco), Dra. Irene Castellón (Universitat de Barcelona) and Dr. Roberto Navigli (Sapienza University of Rome). The thesis was graded Cum Laude.

Keywords: Knowledge acquisition, lexical acquisition, resource evaluation, WordNet, word sense disambiguation

1 Introduction

The use of large-scale semantic resources, such as WordNet, has become a usual, often necessary, practice for most current NLP systems. Princeton WordNet (WN) is by far the most widely-used semantic resource in NLP.

However, even manually, the construction of large-scale semantic repositories for broad-coverage NLP is not a trivial task. It is quite difficult to acquire and consistently integrate large amounts of knowledge into an existing resource. The construction of large and rich knowledge bases for broad-coverage semantic processing takes a great deal of expensive manual effort and incurs large development costs. It involves large research groups for long periods of development. For instance, hundreds of person-years have been invested in the development of wordnets for various languages. In the case of the English Word-

Net, in more than ten years of manual construction (from 1995 to 2006, that is, from version 1.5 to 3.0), WordNet grew from 103,445 to 235,402 semantic relations¹, which represents a growth of around one thousand new relations per month. However, in 2008, the Princeton group released a new resource with 458,825 word forms from the WordNet definitions manually linked to its appropriate WordNet sense².

Furthermore, the relevant knowledge changes across domains and cultures and it has to be steadily kept up to date. New knowledge emerges day by day everywhere and has to be combined with the existing knowledge. For these reasons, knowledge acquisition is still a highly active area of research since the existing knowledge repositories do not

¹Symmetric relations are counted only once.

²<http://wordnet.princeton.edu/glossstag.shtml>

seem to be rich enough to support advanced concept-based NLP applications directly. It seems that such applications require more detailed general-purpose (and also domain-specific) semantic knowledge, which have to be built by automatic means to keep development cost and time inside affordable limits. Obviously, this fact has severely hampered the state-of-the-art of advanced NLP applications.

Thus, the automatic acquisition of the necessary knowledge from available resources, such as naturally occurring text, is one of the most challenging tasks in NLP since it requires some *knowledge understanding* capabilities, which is our final goal. This vicious circle is known as the *acquisition bottleneck*. The intrinsic cycling nature of the problem also suggests a cycling approach for solving it, with incremental iterations of *acquisition-identification* stages. Ideally, the process would start with a *minimal* knowledge base and the relevant resources containing the *implicit* knowledge to be acquired. Then, the automatic acquisition process might produce new content that should be *identified* with respect to the existing knowledge base. This identification process is necessary in order to facilitate the integration of the new knowledge into the existing one, to form a comprehensive and computationally useful knowledge base. Arguably, although these sub-tasks are undeniably difficult, combining them might simplify both.

Figure 1 shows the senses of *party* in WordNet 3.0³. From left to right the figure shows the senses, the total number of explicit semantic relations encoded for each synset, the new semantic relations gathered from the semantically tagged WordNet definitions⁴ and the gloss. Consider the subtle distinctions among some of them. The first three senses are groups of people and the fourth refers to an entertaining event. Obviously, these senses are defining different aspects of related concepts. This is a major drawback when trying to acquire specific knowledge for each sense.

Hopefully, the semantic relations encoded for each sense can help its proper characterization. For instance, Figure 2 shows some of

Sense	#rel.	#gloss	Gloss
party ¹ _n	36	114	an organization to gain political power
party ² _n	18	27	a group of people gathered together for pleasure
party ³ _n	9	41	a band of people associated temporarily in some activity
party ⁴ _n	13	38	an occasion on which people can assemble for social interaction and entertainment
party ⁵ _n	3	87	a person involved in legal proceedings:

Figure 1: Number of relations for party_n in WordNet 3.0

the related concepts encoded in WordNet⁵. Additionally, this table also presents some of the relations captured by KnowNet (KN), a very large lexical knowledge base which has been derived during our research.

Sense	relation	Sense
party ¹ _n	hypernym member-holonym hyponym rgloss related-to	organization ¹ _n , organisation ¹ _n political_system ¹ _n , form_of_government ¹ _n American_Labour_Party ¹ _n machine ⁵ _n , political_machine ¹ _n election¹_n, political¹_a, vote¹_v, elect¹_v
party ² _n	hypernym hyponym hyponym hyponym related-form rgloss related-to	social_gathering ¹ _n , social_affair ¹ _n shindig ¹ _n , shindy ¹ _n dinner ¹ _n , dinner_party ¹ _n wedding ³ _n , wedding_party ¹ _n party ¹ _v carouse ¹ _n carousal ¹ _n bender ² _n toot ² _n , booze-up ¹ _n invitation¹_n, ceremonial¹_a, cocktail¹_n, farewell²_n
party ³ _n	hypernym hyponym rgloss related-to	set ⁵ _n , circle ² _n , band ¹ _n , lot ¹ _n rescue_party ¹ _n fairly ² _r fair ² _r evenhandedly ¹ _r camp⁴_n, landing²_n, stretcher³_n, Olympiad²_n
party ⁴ _n	hypernym hyponym hyponym related-form rgloss related-to	affair ³ _n , occasion ² _n , social_occasion ¹ _n , function ⁶ _n , social_function ¹ _n birthday_party ¹ _n cocktail_party ¹ _n party ¹ _v party-game ¹ _n nuptials¹_n, prom¹_n, reception²_n, gift¹_n
party ⁵ _n	hypernym hyponym domain rgloss related-to	person ¹ _n , individual ¹ _n , someone ¹ _n , somebody ² _n , mortal ¹ _n , soul ² _n assignee ¹ _n law ¹ _n , jurisprudence ² _n submission ⁵ _n accountancy¹_n, appearance³_n, attendance¹_n, court¹_n

Figure 2: Some relations for party_n in WordNet 3.0 and KnowNet(in bold)

1.1 Research goals

The main goal of the research presented in this thesis is to devise new methods and tools

³word^{num}_{pos}, where pos is the part-of-speech (n for nouns, v for verbs, a for adjectives and r for adverbs)

⁴That is, the number of glosses that include that particular sense annotated in its definition

⁵rgloss stands for reverse gloss. That is, the corresponding sense of party appears in its gloss. These relations are gathered from the manually sense-disambiguated glosses

for creating automatically new semantic relations between WordNet senses. That is, to accurately increase by automatic means the knowledge represented in WordNet.

In particular, our research requires the construction of new methods and tools for:

1. Acquiring relevant words from general or domain corpora for an specific WordNet word sense.
2. Identifying the *implicit* word senses of the acquired relevant words with respect to an *existing* knowledge base (in particular, WordNet).
3. Evaluating empirically the quality of the resulting *new* semantic relations in a controlled multilingual evaluation framework.

2 Thesis overview

The thesis is organised in seven chapters:

- **Chapter 1: Introduction**

This chapter presents an overview of the thesis. It revises the motivation and presents the main contributions of the thesis to the state-of-the-art.

- **Chapter 2: State of the Art**

This chapter reviews the state of the art. It revises the use of wide-coverage *semantic resources* in different NLP tasks. Furthermore, it presents the main methodologies, approaches and techniques used for *building large-scale knowledge resources* in general, manually and automatically. Finally, it overviews the existing *evaluation frameworks* used in the research field to assess the quality of the acquired knowledge.

- **Chapter 3: Knowledge Acquisition Method**

This chapter describes the knowledge acquisition architecture developed in this research.

- **Chapter 4: Acquisition of topic signatures**

This chapter reviews the different methods applied to acquire automatically topic signatures as well as the methodology for evaluating their quality.

- **Chapter 5: KnowNet**

This chapter depicts the KnowNet building process and its grounding Word

Sense Disambiguation algorithm, used to obtain word-sense relations from topic signatures acquired from general corpora.

- **Chapter 6: deepKnowNet**

This chapter explores a new method for building KnowNets, named deepKnowNets. Basically, instead of a Word Sense Disambiguation algorithm, the method exploits a graph-based similarity measure to rerank the topic signatures.

- **Chapter 7: Concluding remarks and future directions**

This chapter draws the main conclusions of this thesis and outlines some further steps to follow.

3 Main Contributions

The knowledge acquisition bottleneck problem is particularly acute for open domain (and also domain specific) semantic processing. However, we acquired by fully automatic means highly connected knowledge bases, increasing the total number of semantic relations from less than one million (the current number of available relations in WordNet) to millions of new and accurate semantic relations between WordNet senses. The different versions of KnowNet seem to be a major step towards the autonomous acquisition of knowledge from text, since they are several times larger than the available knowledge resources which encode relations between WordNet senses, and the knowledge they contain outperforms any other resources when they are empirically evaluated in a common framework.

Firstly, in order to acquire relevant semantic relations from large text collections corresponding to general or particular domains, we apply several methodologies and settings to automatically acquire **topic signatures** (TS) (Cuadros, Padró, and Rigau, 2005; Cuadros, Padró, and Rigau, 2006). Originally, topic signatures were used to describe a set of words related to the same topic or domain, but in our case, the topic is a WordNet sense⁶. Thus, topic signatures are sets of words related to that particular WordNet sense. We

⁶The name of Topic Signature, instead of the more appropriate *concept signature*, *word sense signature* or *synset signature*, is maintained for consistency with the literature

use in this research the original topic signatures acquired from the web ⁷ together with new sets of automatically acquired topic signatures which result in new acquisition methods, new tools and different resources, including different types of corpora. All these topic signatures are compared in a common framework together with existing knowledge bases. **ExRetriever**⁸ (Cuadros et al., 2004; Cuadros et al., 2005) is used for the automatic acquisition of examples of particular WordNet word senses. It is a tool to automatically extract a subcorpus of text examples from a large corpus (for instance, BNC, SemCor or the Web).

Secondly, in order to identify the *implicit* semantic relations encoded by a Topic Signature (the sets of words related to a particular WordNet sense) with respect an *existing* knowledge base (in this case, also WordNet), we apply a graph-based Word Sense Disambiguation (WSD) algorithm, SSI-Dijkstra (Cuadros and Rigau, 2008b), also developed in the framework of this thesis. SSI-Dijkstra is based on the Structural Semantic Interconnections (SSI) algorithm. The method uses a robust and accurate knowledge-based Word Sense Disambiguation algorithm to assign the most appropriate senses to semantically related words associated to a particular WordNet sense. In that way, the method identifies a particular WordNet sense for each word in the Topic Signature, converting the original list of concept-to-word relations into a list of concept-to-concept relations.

Thirdly, a variant of SSI-Dijkstra has been applied in a task to integrate a very large domain thesaurus with millions of Species into WordNet (Toral et al., 2010; Cuadros et al., 2010). The process disambiguate every taxonomy of species in several languages.

Finally, the full list of new concept-to-concept relations between WordNet senses forms new knowledge bases, which we call **KnowNet**⁹ (Cuadros and Rigau, 2008b) and **deepKnowNet**. Different sets of new KnowNets are empirically evaluated in different evaluation frameworks (Cuadros and Rigau, 2008b; Cuadros and Rigau, 2008c; Cuadros and Rigau, 2008a; Cuadros and Rigau, 2008d; Agirre et al., 2010).

References

- Agirre, Eneko, Montse Cuadros, German Rigau, and Aitor Soria. 2010. Exploring Knowledge Bases for Similarity. In *Proceedings of LREC 2010*. ISBN: 2-9517408-6-7. Pages 373–377.
- Cuadros, Montse, Jordi Atserias, Mauro Castillo, and German Rigau. 2004. Automatic Acquisition of Sense Examples Using Exretriever. In *Proceedings of IBERAMIA Workshop on Lexical Resources and The Web for Word Sense Disambiguation*, pages 97–104, November.
- Cuadros, Montse, Jordi Atserias, Mauro Castillo, and German Rigau. 2005. The MEANING approach for automatic acquisition of sense examples. *MEANING Workshop*, February.
- Cuadros, Montse, Egoitz Laparra, German Rigau, Piek Vossen, and Wauter Bosma. 2010. Integrating a large domain ontology of species into WordNet. In *Proceedings of LREC 2010*, La Valleta, Malta.
- Cuadros, Montse, Lluís Padró, and German Rigau. 2005. Comparing Methods for Automatic Acquisition of Topic Signatures. In *Proceedings of the International Conference on Recent Advances on Natural Language Processing (RANLP'05)*, September.
- Cuadros, Montse, Lluís Padró, and German Rigau. 2006. An Empirical Study for Automatic Acquisition of Topic Signatures. In *Proceedings of Third International WordNet Conference (GWC 06)*, pages 51–59, Jeju Island (Korea), January. ISBN 80-210-3915-9.
- Cuadros, Montse and German Rigau. 2008a. *Bases de Conocimiento Multilíngües para el Procesamiento Semántico a Gran Escala*. Procesamiento del Lenguaje Natural, Vol. 40, 35–42.
- Cuadros, Montse and German Rigau. 2008b. KnowNet: A proposal for building highly connected and dense knowledge bases from the web. In *First Symposium on Semantics in Systems for Text Processing, STEP'08.*, Venice, Italy, September.
- Cuadros, Montse and German Rigau. 2008c. KnowNet: using Topic Signatures acquired from the web for building automatically highly dense knowledge bases. In *Proceedings of 22nd International Conference on Computational Linguistics (COLING'08)*, Manchester, UK, August.
- Cuadros, Montse and German Rigau. 2008d. Multilingual Evaluation of KnowNet. *Procesamiento del Lenguaje Natural*, 41.
- Toral, Antonio, Monica Monachini, Claudia Soria, Montse Cuadros, German Rigau, Wauter Bosma, and Piek Vossen. 2010. Linking a domain thesaurus to WordNet and conversion to WordNet-LMF. In *Proceedings of ICGL 2010*, Hong Kong.

⁷<http://ixa.si.ehu.es/Ixa/resources/sensecorpus>

⁸<http://www.lsi.upc.edu/nlp/meaning/downloads.html>

⁹<http://adimen.si.ehu.es/web/KnowNet>