Dear Author,

Here are the proofs of your article.

- You can submit your corrections **online**, via **e-mail** or by **fax**.

- For **online** submission please insert your corrections in the online correction form. Always indicate the line number to which the correction refers.

- You can also insert your corrections in the proof PDF and **email** the annotated PDF.

- For fax submission, please ensure that your corrections are clearly legible. Use a fine black pen and write the correction in the margin, not too close to the edge of the page.

- Remember to note the **journal title**, **article number**, and **your name** when sending your response via e-mail or fax.

- **Check** the metadata sheet to make sure that the header information, especially author names and the corresponding affiliations are correctly shown.

- **Check** the questions that may have arisen during copy editing and insert your answers/corrections.

- **Check** that the text is complete and that all figures, tables and their legends are included. Also check the accuracy of special characters, equations, and electronic supplementary material if applicable. If necessary refer to the *Edited manuscript*.

- The publication of inaccurate data such as dosages and units can have serious consequences. Please take particular care that all such details are correct.

- Please **do not** make changes that involve only matters of style. We have generally introduced forms that follow the journal's style.
  Substantial changes in content, e.g., new results, corrected values, title and authorship are not allowed without the approval of the responsible editor. In such a case, please contact the Editorial Office and return his/her consent together with the proof.

- If we do not receive your corrections **within 48 hours**, we will send you a reminder.

- Your article will be published **Online First** approximately one week after receipt of your corrected proofs. This is the **official first publication** citable with the DOI. **Further changes are, therefore, not possible.**

- The **printed version** will follow in a forthcoming issue.

**Please note**

After online publication, subscribers (personal/institutional) to this journal will have access to the complete article via the DOI using the URL: http://dx.doi.org/[DOI].
If you would like to know when your article has been published online, take advantage of our free alert service. For registration and further information go to: http://www.springerlink.com.

Due to the electronic nature of the procedure, the manuscript and the original figures will only be returned to you on special request. When you return your corrections, please inform us if you would like to have these documents returned.

# Metadata of the article that will be visualized in OnlineFirst

| | |
|---|---|
| ArticleTitle | Challenges for a multilingual wordnet |
| Article Sub-Title | |
| Article CopyRight | Springer Science+Business Media B.V. <br> (This will be the copyright line in the final PDF) |
| Journal Name | Language Resources and Evaluation |

| Corresponding Author | Family Name | **Fellbaum** |
|---|---|---|
| | Particle | |
| | Given Name | **Christiane** |
| | Suffix | |
| | Division | Department of Computer Science |
| | Organization | Princeton University |
| | Address | Princeton, NJ, 08540, USA |
| | Email | fellbaum@princeton.edu |

| Author | Family Name | **Vossen** |
|---|---|---|
| | Particle | |
| | Given Name | **Piek** |
| | Suffix | |
| | Division | Faculty of Arts |
| | Organization | VU University of Amsterdam |
| | Address | Amsterdam, 1081 HV, The Netherlands |
| | Email | |

| Abstract | Wordnets have been created in many languages, revealing both their lexical commonalities and diversity. The next challenge is to make multilingual wordnets fully interoperable. The EuroWordNet experience revealed the shortcomings of an interlingua based on a natural language. Instead, we propose a model based on the division of the lexicon and a language-independent, formal ontology that serves as the hub interlinking the language-specific lexicons. The ontology avoids the idiosyncracies of the lexicon and furthermore allows formal reasoning about the concepts it contains. We address the division of labor between ontology and lexicon. Finally, we illustrate our model in the context of a domain-specific multilingual information system based on a central ontology and interconnected wordnets in seven languages. |
|---|---|

| Keywords (separated by '-') | Multilingual wordnets - Formal ontology - Information system |
|---|---|
| Footnote Information | |

Springer
the language of science

# Author Query Form

## Please ensure you fill out your response to the queries raised below and return this form along with your corrections

Dear Author

During the process of typesetting your article, the following queries have arisen. Please check your typeset proof carefully against the queries listed below and mark the necessary changes either directly on the proof/online grid or in the 'Author's response' area provided below

| Query | Details required | Author's response |
|---|---|---|
| 1. | Please confirm the section headings are correctly identified. | |
| 2. | Please check and confirm the volume number and issue number is correctly identified for the reference Gangemi et al. (2003). | |
| 3. | Please check and confirm the volume number, issue number and page range are correctly identified for the references Miller (1990), Tufis (2004). | |
| 4. | Please check and confirm the page range is correctly identified for the reference Miller (1995). | |

ORIGINAL PAPER

# Challenges for a multilingual wordnet

**Christiane Fellbaum · Piek Vossen**

**Abstract** Wordnets have been created in many languages, revealing both their lexical commonalities and diversity. The next challenge is to make multilingual wordnets fully interoperable. The EuroWordNet experience revealed the short-comings of an interlingua based on a natural language. Instead, we propose a model based on the division of the lexicon and a language-independent, formal ontology that serves as the hub interlinking the language-specific lexicons. The ontology avoids the idiosyncracies of the lexicon and furthermore allows formal reasoning about the concepts it contains. We address the division of labor between ontology and lexicon. Finally, we illustrate our model in the context of a domain-specific multilingual information system based on a central ontology and interconnected wordnets in seven languages.

**Keywords** Multilingual wordnets · Formal ontology · Information system

## 1 Introduction

Digital lexical resources can store lexicons of potentially unlimited size in ways that enable flexible representations and searches. Mapping the lexical inventory of a language into a semantic network has proved to be useful for many natural language processing applications, and WordNet-style lexicography has been applied to build

A1   C. Fellbaum (✉)
A2   Department of Computer Science, Princeton University, Princeton, NJ 08540, USA
A3   e-mail: fellbaum@princeton.edu

A4   P. Vossen
A5   Faculty of Arts, VU University of Amsterdam, 1081 HV Amsterdam, The Netherlands

26 resources in many languages.[1] The challenge we face now is to interconnect them so
27 as to create one multilingual database. To reflect intra- and inter-lingual meaning,
28 we argue for the connection of lexical databases to a shared ontology, which
29 requires us to resolve some fundamental linguistic and ontological questions. We
30 address these in the context of an ongoing project that represents a first step in the
31 creation of a global wordnet system.

## 2 The scope of a lexicon

33 Lexical databases do not need to obey constraints on their size, and no well-defined
34 guidelines exist for what is to be included in the lexicon. Lexicons are idiosyncratic;
35 they do not systematically label concepts, and the idiosyncracies are often revealed
36 in crosslinguistic differences. Lexicons are moreover redundant, often assigning
37 several word forms (synonyms) to a single concept. Finally, lexicons are open-
38 ended, often extended into terminology and domain-specific vocabulary.
39  Because inclusion in the lexicon cannot be determined by well-defined rules, its
40 boundaries are fuzzy. Moreover, the lexical status of many phrases and chunks is
41 uncertain, raising the question as to what constitutes a lexeme deserving of a
42 legitimate entry in the databases.
43  While even linguistically naive speakers have an intuitive notion of "word," there
44 exists no hard definition. A possible orthographic definition would state that strings
45 of letters with an empty space on either side are words. While this would cover
46 words such as *road, eat*, and *heavy*, it would wrongly leave out multiword units like
47 *lightning rod, find out, word of mouth*, and *spill the beans* that constitute semantic
48 and lexical units.[2] A first rule of thumb might state that a lexical unit will merit
49 inclusion in a database when it serves to denote an identifiable concept. But this
50 criterion is less than straightforward, especially when applied to multi-word units.

51 2.1 Compositionality, currency, salience, linguistic authority

52 For NLP applications, including multi-word phrases and sentence chunks besides
53 single words may present clear processing advantages. However, even so-called
54 "fixed expressions" are subject to lexical variation and internal modification (e.g.,
55 Fellbaum 2007). The lexical status of multiword units is usually determined on the
56 basis of the compositionality criterion. If the meaning of the whole is the sum of the
57 meaning of its parts, there is no reason to consider the multiword string a separate
58 entity. Thus, fully compositional phrases like *coniferous forest* and *subtropical*
59 *coniferous forest* should probably not be considered as separate fixed lexical items.
60 By contrast, a non-compositional term like *ecological footprint* cannot be readily
61 analyzed by speakers and must be listed in the lexicon. Additional factors, not based

---

1FL01 [1] We will reserve the legally registered name "WordNet" for the Princeton WordNet and use "wordnet"
1FL02 as a generic term to denote semantic networks inspired by the Princeton WordNet.

2FL01 [2] In languages whose writing systems do not separate lexical units, the notion of word is of course
2FL02 divorced from that of a graphemic unit.

62 on linguistic economy, might override the compositionality criterion. *Currency,*
63 *salience,* and *speaker authority* are three such interacting factors.

64    Currency is the extent to which a speaker community avails itself of a word or
65 phrase that becomes (often temporarily) salient through frequent use. While
66 frequency and shared cultural background determine the currency of a word or
67 phrase, the authority of a speaker or a subgroup of speakers within a language
68 community may have an effect on a word's use as well. Thus, popular media
69 exercise a significant influence on the words that are circulating within a speaker
70 community; frequency counts for a given lexeme vary over time, as the
71 newsworthiness of stories and topics grows and diminishes. Social groups determine
72 acceptance and linguistic change, as studies of youth language have shown.

73    Within a specific domain, a multiword term may be particular salient, as reflected
74 in its frequency or its function as a topic of numerous documents. Over time, such
75 compounds may acquire the status of fixed, ready-made expressions and become a
76 part of the lexicon of the language community. Compound terms become
77 established in a language community when their creators and initial users have a
78 social standing that bestows on them a kind of "linguistic authority." This
79 phenomenon can be seen in the areas of science and technology (*mouse potato*),
80 popular entertainment and commercial branding (*e-inkReader*), where people
81 introduce new terms often with the wish of adding them, along with a new concept,
82 to the lexicon.

83    An example of currency, salience and perhaps linguistic authority is the Dutch
84 noun compound *arbeidstijdverkorting*. Although its members, *arbeid* ("work"), *tijd*
85 ("time"), and and *verkorting* ("reduction") suggest a straightforward compositional
86 meaning, this compound in fact denotes more than the mere sum of its members: a
87 specific social arrangement dating to the 1980s intended to decrease unemployment;
88 work hours and wages were reduced so that additional workers could be hired and
89 new jobs could be created.

90 **3 WordNet, EuroWordNet, global wordnet**

91 Digital lexicography resulted in the abandonment of orthography as an organizing
92 principle for dictionaries. Representing the lexicon as a semantic network has
93 proved particularly useful for NLP applications, and WordNet-style resources were
94 built for many languages. We briefly review the principles of wordnet lexicography
95 and the development of multilingual wordnets.

96    The Princeton WordNet (Miller 1990, 1995; Fellbaum 1998) is a manually
97 constructed large-scale lexical database for English. WordNet's original motivation
98 was to test the feasibility of a model of human semantic memory that sought to
99 explain principles of storage and retrieval of words and concepts. This model
100 proposed a largely hierarchical organization of concepts expressed by nouns, events
101 (encoded by verbs) and properties (expressed by adjectives). The WordNet
102 experiment tried to determine whether the bulk of the lexicon of a language could
103 indeed be represented in a semantic network by means of a handful of relations,
104 inspired by the model of human memory.

105 An important semantic relation is that between words sharing the same
106 denotational meaning, synonymy. WordNet groups synonymous words into
107 "synsets," the building blocks, or nodes, of the semantic network. Synsets are
108 interlinked by means of semantic relations, such as hyponymy (the super-
109 subordinate relation that holds between words like *building* and *garage*), meronymy
110 (the part-whole relation that hold between words like *toe* and *foot*), antonymy (the
111 relation between contrasting concepts, such as *expensive* and *cheap*) and troponymy
112 (a "manner" relation that links verbs like *prance* and *walk*). In addition to relations
113 among synsets, WordNet also includes lexical relations among specific synset
114 members—morphologically and semantically related form-meanings pairs such as
115 *direct, director*, and *direction*, etc. (Fellbaum and Miller 2003).

116 WordNet defines membership in a synset as denotational equivalence and
117 substitutability in some, though not all, contexts. But in fact synsets include
118 members that differ along many dimensions, and some are more similar to one
119 another than others. A more subtle representation could label the many ways synset
120 members are related to one another, such as in terms of dialectal variations as in the
121 synsets {*grinder, hero, hoagie, poor boy, submarine*} or register, ranging from
122 formal to taboo words whose use is restricted to particular settings.

123 Although it was not motivated by Natural Language Processing research,
124 WordNet's extensive coverage, digital format, and above all its graph structure
125 make it useful for automatic language processing. When WordNet was widely
126 embraced by the English NLP community, wordnets construction began in other
127 languages.

128 Within the EuroWordNet project (Vossen 1998), lexical databases modeled on
129 the Princeton WordNet were constructed for eight languages. EuroWordNet
130 contributed several fundamental innovations to the wordnet design that have since
131 been adopted by dozens of additional wordnets. One is the definition of a set of Base
132 Concepts, which are characterized by many links to other synsets in wordnets and
133 which are assumed to be universally encoded. Second, to increase the connectivity
134 among synsets, a number of new relations were defined, in particular cross-part-of-
135 speech relations. All relations were marked with a feature value indicating the
136 combinations of relations (conjunctive or disjunctive) and their directionality.
137 Relations may be language-specific rather than apply to all wordnets.

138 Following EuroWordNet, wordnets were developed for a number of languages
139 around the world. Besides individual efforts, there are wordnets for entire
140 geographic regions, such as BalkaNet (Tufis 2004), African Wordnet (Moropa
141 et al. 2007), Asian wordnet (Robkop et al. 2010) and the Indian wordnets (Sinha
142 et al. 2006). Currently, wordnets exist for some sixty genetically and typologically
143 distinct languages (cf. www.globalwordnet.org).

144 Many wordnets are created independently and their coverage and design is not
145 uniform. The challenge is to create a framework that accommodates an ever-
146 increasing diversity of languages without shoehorning them into a pattern developed
147 originally for English only. Fellbaum and Vossen (2007) and Vossen and Fellbaum
148 (2009) present an outline of the Global WordNet Grid, a proposed system designed
149 to accommodate the lexicons of all languages and connect them via a language-
150 independent ontology.

| Journal : **10579** | Dispatch : **30-4-2012** | Pages : **14** |
| Article No. : **9186** | □ LE | □ TYPESET |
| MS Code : **LRE881** | ☑ CP | ☑ DISK |

Challenges for a multilingual wordnet

151 3.1 Language-specific relations

152 Synonymy, at least under the somewhat loose definition that WordNet adopts,
153 appears to be a universal phenomenon. And we have not yet encountered a language
154 whose lexicon cannot be organized at least partly by means of hyponymy,
155 meronymy, and antonymy. But some other semantic distinctions are lexicalized in a
156 subset of the world's languages only. This poses challenges for representing the
157 lexemes in a language-independent, universally valid ontology.

158 *3.1.1 Gender lexicalizations*

159 Consider male and female profession nouns like *actor-actress*. While English
160 does not distinguishes gender systematically and nouns like *teacher, lover, friend,*
161 etc. are underspecified with respect to gender, languages like German and Dutch
162 systematically and regularly encode it. The female form is usually derived in a
163 productive fashion by means of a suffix. Because this process is not shared by all
164 languages, only those that mark the distinction will link the noun pairs via a
165 gender relation in their lexicons. The members of such pairs connect to the
166 corresponding classes in the ontology ("male" or "female"); English words are
167 linked to both.

168 *3.1.2 Verbal aspect*

169 Verbal aspect is distinguished and encoded differently across languages. Languages
170 including English and German can differentiate activities from accomplishments by
171 adding a particle to a simplex verb, as in the English pairs *eat up* and *read through*
172 (German *aufessen* and *fertiglesen*). Perfectivity is not obligatorily marked, and a
173 sentence like *Peter read a magazine* is underspecified as to whether or not Peter
174 read the entire magazine, front to back.
175 Other languages regularly encode semantic distinctions by means of affixes.
176 For example, Slavic languages systematically distinguish between the imperfec-
177 tive, unprefixed and the perfective, prefixed members of a verb pair. Czech has
178 1,000s of such verb pairs, where one member derived via regular and productive
179 morphology. Do aspectual distinctions belong into the lexicon or the ontology?
180 The fact that they are not universally marked (Romance, for example, uses
181 different conjugational endings but no lexical encoding) might argue for a
182 relation among aspectually related verb pairs in the lexicons of German, English,
183 Czech, etc. (Pala et al. 2008). All verb forms related to the same base form
184 would be linked to one event in the ontology. However, limiting the encoding of
185 verbal aspect to the lexicon and excluding it from the ontology will constrain the
186 reasoning power of the ontology (e.g., if completive *eat up* is not distinguished
187 from progressive *eat*, a system cannot draw any conclusions as to whether the
188 food has been completely consumed or not).

### 3.1.3 Event perspective

Some events involving multiple participants can be expressed with different verbs that profile different participants in the event. For example, converse pairs like *buy* and *sell* express the actions of different participants in the same sale event. FrameNet (Ruppenhofer et al. 2002) captures this difference by referring to distinct Frame Elements—Buyer and Seller—of a single Frame.

While the verbs and the corresponding Agent nouns (*buyer, seller*) each merit their own lexical entries, in the ontology they can be represented so as to reflect different perspectives on the same event. Converse and reciprocal events may be encoded very differently across languages. While English labels the two sides of a sale event with distinct word forms (*buy*, *sell*), others, like German, distinguish them by means of a morpheme (*kaufen* vs. *verkaufen*). And whereas English encodes the difference between the activities of a teacher and a student in two different verbs, *teach* and *learn*, French uses the same verb, *apprendre*, and encodes the distinction syntactically.

Russian has two different verbs corresponding to English *marry*, depending on whether the grammatical subject refers to the bride or the groom. In such cases, the lexicons need only refer to the event entry in the ontology (*sale*, *marriage*, etc.) and implement equivalence mappings between the terms and lexical entities, leaving the linguistic encoding of distinct verbs and roles to the lexicons of each language.

Crosslinguistic lexicalization patterns show the need for a broader, language-independent treatment that can accommodate all variations on the language level but unifies them on the conceptual, ontological level.

## 4 Natural language interlingua

Because the lexicons of different languages do not all label the same concepts, a simple mapping from English to the target languages and across the EuroWordNet languages is ruled out in many cases. To interconnect the wordnets, EuroWordNet linked the synsets of each language via an "equivalence relation", to an interlingual index, or ILI. The ILI permits the mapping of equivalent synsets across all languages connected to the ILI, and thus allows not only for straightforward translations but also for the comparison of the lexicons of different languages both in terms of coverage, relations, and overall lexicalization patterns.

Initially, the ILI was populated with the synsets from the Princeton WordNet, which provided large coverage and was accessible to speakers of the EuroWordNet languages, enabling them to judge semantic equivalence.

EuroWordNet revealed the problems that arise when a natural language becomes the hub connecting the lexicons of other languages. The first concerns *coverage.* No two languages have completely overlapping lexicons. For many concepts, one language may have one or more lexical labels while another language has none. An ILI tied to one specific language clearly reflects only the inventory of the language it is based on, and gaps show up when lexicons of different languages are mapped to it. Using a natural language as the interlingua also may bias the coverage and

| Journal : **10579** | Dispatch : **30-4-2012** | Pages : **14** |
| Article No. : **9186** | ☐ LE | ☐ TYPESET |
| MS Code : **LRE881** | ☑ CP | ☑ DISK |

Challenges for a multilingual wordnet

231 representation of the wordnets of other languages. Interestingly, those EuroWordNet
232 languages that translated the English WordNet (using the "Expand" method)
233 constructed different wordnets from those that started independently and later
234 mapped onto the ILI (the "Merge" method).

235 More serious is the question of *equivalence*. The semantic space covered by a
236 word in one language often overlaps only partially with that covered by a similar
237 word in another language, making for less than perfect mappings. An apparently
238 good crosslinguistic match may turn out not to be one when one considers different
239 contexts and social settings. This is the case for connotational differences, tied to
240 specific usages of the words. Second, the mappings among the words and synsets in
241 the ILI may appear to be appropriate on the word level, but there may be a
242 difference in their position within their respective local networks. Such a mismatch
243 necessarily reflects a meaning difference, since in a semantic network the meaning
244 of a node is by definition given in terms of its relations to other nodes. For example,
245 the fact that Dutch lacks a word for "container" does not mean that *bag, box, bottle*
246 etc. do not form a natural category in Dutch, as they do in English by virtue of being
247 children of *container*.

248 Finally, although WordNet borrows relations like hyponymy and meronymy
249 from ontology, it does not encode the lexicon with such relations in ways that reflect
250 clean *ontological methodology*. As Guarino and Welty (2002a, b) and Gangemi
251 et al. (2003), among others, point out, WordNet's hyponymy relation includes
252 multiple, distinct relations. Earlier versions conflated types, instances, and roles.
253 Thus, *Bill Clinton* was "a type of" *President*, just as *desk* was "a type of" *table*. A
254 later version drew the distinction between Types and Instances, so that proper
255 names referring to people, products, countries, mountains, stars, etc. are now all
256 Instances (Miller and Hristea 2006) and only common nouns can be Types.
257 However, Roles are not presently distinguished from Types, so that *president* and
258 *professor* continue to be represented as "types of" *person* (cf. Sect. 5 for further
259 discussion).

## 5 From interlingua to ontology

261 Arguably, using a language-independent interlingua as the hub that connects
262 language-specific lexicons is a better approach to mapping lexicons than a direct
263 mapping. But the interlingua must be able to represent concepts expressed by words
264 in a way that is not biased towards any language or any word-specific linguistic
265 properties at all. The division between words and concepts is reflected in that
266 between the lexicon and ontology.

267 The use of Princeton WordNet as the interlingua in EuroWordNet blurred this
268 distinction, and the KYOTO project described in Sect. 6 aim to restore it by
269 assigning words on the one hand to wordnet-like structured lexicons and by
270 relegating concepts to ontology.

271 Lexicons–mappings of labels (words, or lexemes) to concepts (mental represen-
272 tations of entities)—are natural, not products of human reasoning or reflection. They
273 have an internal structure, which is revealed by (often productive) lexicalization

274 patterns and distinct linguistic properties for lexical subclasses (e.g., Levin 1993).
275 But lexicons have many idiosyncrasies, such as seemingly unmotivated, "acciden-
276 tal" gaps. Lexicons also show that languages tend to have several labels for given
277 concept (synonymy), though the words may not all be fully equivalent. While the
278 lexicons of all languages may share a core concept-word mapping inventory,
279 language-specific idiosyncracies abound.

280 WordNet is often called a lexical ontology because it records lexicalized
281 categories and connects them by means of relations familiar from formal ontology.
282 However it differs in significant ways from a formal ontology, an artificially
283 constructed design. Ontologies are language-independent; the linguistic labels in
284 their axioms are merely conveniences and are not to be confused with words used in
285 a natural language. Consequently, the mapping from lexicon to ontology is one from
286 word to concept, rather than across words and languages as in the case of the
287 EuroWordNet ILI. Ontology aims to be completely unambiguous about the meaning
288 of its entries, whereas word meanings are typically fuzzy. Moreover, ontological
289 relations do not necessarily reflect speakers' intuitions about relations among words.

290 Because each of its entries is unique, clearly defined and distinguished from
291 every other entry. Ontology is preferable over a language-specific lexicon as the hub
292 connecting wordnets of different languages, as argued by Fellbaum and Vossen
293 (2007), Vossen and Fellbaum (2009), and Pease and Fellbaum (2009). This allows
294 for a clean separation between the lexicons and a language-independent, formal
295 representation of the concepts lexicalized by individual wordnets. Moreover, the
296 burden of expressing relations among words and formal concepts can be shared
297 between the lexicons and the ontology. The SUMO ontology (Niles and Pease 2001;
298 2003) was the first to have been mapped to a number of wordnets and to function as
299 their interlingua.

## 6 Ontology

301 In the context of artificial intelligence (AI) and knowledge engineering, ontology is
302 the explicit, formal specification of a conceptualization (Gruber 1992; 1993). For AI
303 systems, what "exists" is that which can be represented. A formal ontology contains
304 definitions that associate the names of entities in the universe of discourse (e.g.,
305 classes, relations, functions, or other objects) with human-readable text describing
306 what the names mean, and formal axioms that constrain the interpretation and well-
307 formed use of these terms; furthermore, ontology specifies the relations among
308 concepts (see e.g., Gruber 1993).

309 The ontology takes input from the lexicons, but on a "selective"' basis, such that
310 not all lexicalized entities are added to the ontology. While the ontology must be
311 able to encode all concepts that can be expressed in any natural language, it need not
312 provide a linguistic encoding—a label—for all words and expressions.

313 It is desirable that the ontology contain only terms distinguished by essential
314 properties; second, that it be comprehensive and include all distinct concepts that
315 can be represented as Types for all languages; third, that equivalent concepts across
316 languages can be related; fourth, that it allow the definition of all lexicalized

| Journal : **10579** | Dispatch : **30-4-2012** | Pages : **14** |
| Article No. : **9186** | □ LE | □ TYPESET |
| MS Code : **LRE881** | ☑ CP | ☑ DISK |

Challenges for a multilingual wordnet

317 concepts having non-essential properties, and finally, that it be logically valid and
318 allow for inferencing.

319   Guarino and Welty (2002a, b) demonstrated that the WordNet hierarchy, when
320 examined with ontological criteria, can be improved and reduced. Their proposed
321 OntoClean method relies on metaproperties to determine the ontological properties
322 of classes and can be applied to determine the smallest common set of concepts in
323 all languages. The properties of these concepts are *rigidity, essence, dependence* and
324 *unicity*.

325   Guarino and Welty's rigidity criterion is particularly relevant for the consistent
326 distinction between lexicon and ontology, because languages encode many non-
327 rigid concepts. Rigidity distinguishes Types such as *poodle, Newfoundland, German*
328 *shepherd* from Roles like *lapdog* and *herding dog*. Types and Roles are not disjunct:
329 a given entity may be both a Type and and a Role at the same time. While a German
330 shepherd will never be a Newfoundland or a poodle, German shepherds may assume
331 different Roles such as that of a herding dog or a lap dog. Only types of dogs are
332 included in the ontology; if a language lexicalizes a role such as *herding dog*, the
333 type hierarchy of the ontology is not extended, but the word is defined in the
334 ontology and marked as a Role (Vossen et al. 1999).[3]

335   One could include in the ontology all the relations that are found in a semantic
336 network like WordNet. Having done that, the question would be how to express
337 informal linguistic notions with more formal ontological relations. By keeping
338 ontological relation in the formal ontology, and linguistic relations in the lexicon,
339 one can avoid merging two different levels of analysis and yet still capture the
340 information that is needed about both formal concepts and linguistic tokens. An
341 important requirement for the ontology is that it be suitable for automatic reasoning.
342 Therefore, relations in the ontology must be logically consistent and apply strictly.

343   In a lexicon or a semantic network the meaning of a word can be expressed with
344 natural language definitions. Word meanings as represented in a lexicon are subject
345 to human judgment and introspection. By contrast, in ontology it is solely the
346 axioms as formal statements that gives the terms their meaning. Although the
347 axioms borrow words from natural language, the meanings of these terms are
348 independent of their surface forms. One could replace all the term names with
349 arbitrary unique symbols and they would still have the same meaning. This entails
350 that the meaning of the terms can be tested for consistency with an automated
351 theorem prover, rather than the ontologist having to rely completely on human
352 inspection and judgments of word meaning.

353 **7 Case study: KYOTO, a multilingual information system**

354 KYOTO (Knowledge-Yielding Ontologies for Transition-Based Organization), a
355 project funded by the European Union's Seventh Framework (http://www.kyoto-
356 project.eu), represents the first step toward a Global WordNet. KYOTO rests on

---

3FL01   [3] A small number of salient and possibly universally lexicalized roles, including *mother, father, friend*
3FL02   will be included in the type hierarchy.

357 the twin pillars of formal concept representations (ontology) and linguistic
358 representations (lexicons, wordnets), whose division and interrelations allow one
359 to build a domain-specific multilingual wordnet system anchored in a language-
360 independent central ontology. The system is designed to allow easy crosslingual
361 sharing and transfer of information both by automatic systems and by human
362 users without a background in Knowledge Engineering. It enables its users to
363 build crosslinguistic consensus on the meaning and interpretation of language.
364 KYOTO is validated for specific, interlocking domains including biodiversity,
365 climate change and environmental protection (Vossen et al. 2008).

## 7.1 The KYOTO architecture

367 KYOTO uses a three-layered knowledge model that separates (1) multilingual
368 general and domain-specific vocabularies linked to (2) multilingual generic and
369 domain-specific wordnets connected to the English WordNet, and (3) a language-
370 independent, formal central ontology, to which all wordnets are linked. Each layer
371 has an internal semantic structure that allows one to connect specific concepts to
372 more general concepts via explicit explicit mapping relations. The ontology
373 contains rich axioms for modeling processes and qualities.
374   In a first step, human experts identify and specify the locations and sources of
375 domain-relevant documents in different languages. Term extraction from these texts
376 is performed by linguistic miners, so-called term-yielding robots ("tybots"), which
377 identify relevant domain terms and the concepts behind them and relate them to
378 semantic networks (wordnets) in English, Dutch, Spanish, Basque, Italian, Chinese,
379 Japanese. The miners identify possible relations (such as hyponymy) among the
380 members of a phrase or a compound. For example, the miners can suggest that *water*
381 is the polluted entity in the term *water pollution*.
382   A wiki environment allows ontologically "naïve" users to add domain terms in a
383 way that respects important distinctions among concepts, in particular Rigidity. An
384 editor prompts the domain-experts to identify and encode formal constraints and
385 relations among the terms representing entities, processes and states. This results in
386 a computationally tractable domain ontology that is made available to other user
387 communities where cross lingual validation takes place. The domain wordnets and
388 the ontology are harmonized and anchored to general-coverage wordnets and a
389 generic (domain-independent) ontology.

## 7.2 The KYOTO ontology

391 A central question for the system concerns the division of labor between the
392 language-specific lexicons and the ontology (Vossen and Rigau 2010). We outline
393 the criteria for building and distinguishing these two key components of the system.
394   A top-level ontology is defined as well as a middle level ontology that makes it
395 possible to integrate the environmental knowledge of the applied domain. It would

396 be impossible to represent in the wordnets and in the ontology all complex terms
397 found in domain-specific databases and texts, let alone to attempt automatic
398 inferencing over the terms. Therefore, only a subset of the concepts are represented
399 in the domain-specific wordnets and the generic ontology (which contains only rigid
400 entities) while more specific terms are linked to these via subsumption relations. As
401 a result, the ontology is the direct hub for only a subset of the concepts. In addition,
402 KYOTO makes the assumption that the generic wordnets and vocabularies contain
403 mostly rigid types (e.g., *frog*), whereas domain-specific documents with news and
404 event-specific information typically include in addition non-rigid concepts such as
405 *endangered frogs, endemic frogs* and *alien frogs*. KYOTO allows one to distinguish
406 the rigid entities referred to by a substring of such expressions (e.g., *frog*) and to
407 identify their semantic relation to the states and processes expressed by the
408 remaining constituents (e.g., *endangered*).
409     A number of mapping relations relate the expressions referring to states and
410 processes in the generic wordnets to the appropriate entries in the ontology.

411 7.3 Mapping between wordnets and the central ontology

412 The ontology can represent the processes, states and qualities that are relevant for
413 the KYOTO domain. Mappings were created for highly frequent verbs and
414 adjectives in the domain (e.g., *endanger, endemic*) to these processes, states and
415 qualities in order to differentiate between rigid and non-rigid concepts in the
416 wordnets and to be able to match the non-rigid concepts to the relevant
417 processes. As an example, consider the term *migratory bird*. To reflect that this
418 non-rigid term is a hyponym of *bird* but not a proper subclass, the following
419 mapping was created:

420     *wn:migratory bird* sc_domainOf *ont:bird*
421     *wn:migratory bird* sc_playRole *ont:done-by*
422     *wn:migratory bird* sc_participantOf *ont:migration*

423     This mapping indicates, first, that the term is used to refer to instances (but not
424 subclasses) of endurants, where the domain is restricted to birds. In addition, the
425 mapping states that the concept in question participates in the process of migration
426 as a participant (in the role of done-by).
427     The process "migration" is further defined in the ontology, stating that it is an
428 active-change-of-location done-by some endurant, going from a source via a path to
429 some destination. The mapping relations from the wordnet to the ontology need to
430 satisfy the constraints of the ontology, i.e. only roles can be expressed that are
431 compatible with the role-schema of the process in which they participate. The
432 wordnet-to-synset mappings can thus be used to define fairly basic relations relative
433 to the ontology, which represents the full meanings of the terms.
434     These mappings can clarify many subtle meaning differences among closely
435 related concepts across languages. Consider the following examples:

438 {wn:teacher} English                  {wn:meat} English

439 →sc_domainOf **ont:human**           →sc_domainOf **ont:cow, sheep, pig**

440 →sc_playRole ont: **ont:done-by**       →sc_playRole **ont:patient**

441 →sc_participantOf **ont:teach**         →sc_participantOf **ont:eat**

442 {wn:leraar} Dutch // lit. *male teacher*    {wn:名 肉, 食物, 餐 } Chinese

443 →sc_domainOf **ont:man**               →sc_domainOf **ont:animal**

444 →sc_playRole **ont:done-by**         →sc_playRole **ont:patient**

445 →sc_participantOf **ont:teach**         →sc_participantOf **ont:eat**

446 {wn:lerares} Dutch // lit. *female teacher*   {wn: غذاء, لحم, طعام} Arabic

447 →sc_domainOf **ont:woman**          →sc_domainOf **ont:cow, sheep**

448 →sc_playRole **ont:done-by**         →sc_playRole **ont:patient**

449 →sc_participantOf **ont:teach**         →sc_participantOf **ont:eat**

450 On the left, we see mappings for English and Dutch synsets to the role of a
452 *teacher*, where the domain in English is restricted to humans but in Dutch it is
453 differentiated into men and women. On the right, we see representations for edible
454 kinds on meat in English, Chinese and Arabic; note that the domains differ across
455 these languages. The EuroWordNet ILI solution required a mapping from all the
456 non-English synsets to the English ones, blurring often important differences;
457 moreover, it would not allow a flexible representation of non-rigid concepts as in the
458 example above. The solution in KYOTO allows us to keep the differences explicit
459 and at the same time keep the ontology restricted.

460 7.4 Reasoning and inferencing with KYOTO

461 The reasoning and inferencing capabilities of KYOTO incorporate the three-layered
462 knowledge model and the notion of an explicit ontology in which a relevant subset
463 of implications is exported to be inserted into knowledge annotation format
464 representations of text.
465 As the example above shows, classes in the ontology are defined using rich
466 axioms that specify the semantics needed for inferencing: "migration" is represented
467 as an *active-change-of-location done-by* some endurant, going from a source via a
468 path to a destination. At any given point in the ontology development, KYOTO
469 creates an *explicit ontology*, which is a collection of all the implications that apply
470 to a class given the OWL-DL specification of the ontology.
471 Different surface forms like *migratory birds, bird migration, migration of bids,
472 birds that migrate* are subject to the same ontological implications that build on the
473 relation between the migration process and birds and also provide place holders for
474 other elements in the text to map to the source, path and destination. The same holds
475 for processing of text in languages other than English: regardless of the linguistic
476 (morphosyntactic) structure expressing a concept, the ontology provides the same
477 semantic model.

| | Journal : **10579** | Dispatch : **30-4-2012** | Pages : **14** |
|---|---|---|---|
| [logo] | Article No. : **9186** | □ LE | □ TYPESET |
| | MS Code : **LRE881** | ☑ CP | ☑ DISK |

Challenges for a multilingual wordnet

## 8 Summary and conclusion

There are multiple challenges for aligning wordnets for different languages and create a system that allows crosslinguistic mapping and facilitates automatic language processing. The overall design imposes a clear division between the language-specific lexicons (wordnets) and a formal, language-independent ontology that serves as the hub by which to which all wordnets are interconnected. Ontology is constructed according to strict principles, while the lexicons show a variety of idiosyncracies with respect to the linguistic encoding of concepts and lexical patterns. The KYOTO project provides a framework for the division of labor between ontology and lexicons and for the empirical investigation of the kinds of lexical and sublexical information that ontology can efficiently represent.

## References

Fellbaum, C. (Ed.). (1998). *WordNet: An electronic lexical database*. Cambridge, MA: MIT Press.

Fellbaum, C. (Ed.). (2007). *Collocations and idioms: Corpus-based linguistic and lexicographic studies*. Birmingham, UK: Continuum Press.

Fellbaum, C., & Miller, G. A. (2003). Morphosemantic links in WordNet. *Traitement automatique de langue, 44*(2), 69–80.

Fellbaum, C., & Vossen, P. (2007). Connecting the universal to the specific. In T. Ishida, S. R. Fussell & P. T. J. M. Vossen (Eds.), *Intercultural collaboration; First international workshop* (Vol. 4568, pp. 1–16). Lecture Notes in Computer Science, Springer, New York.

Gangemi, A., Guarino, N., Masolo, C., & Oltramari, A. (2003). Sweetening WordNet with DOLCE. *AI Magazine, 57*(1), 13–24.

Gruber, T. R. (1992). A translation approach to portable ontologies. *Knowledge Acquisition, 5*(2), 199–220.

Gruber, T. R. (1993). A translation approach to portable ontologies. *Knowledge Acquisition, 5*, 199–220.

Guarino, N., & Welty, C. (2002a). Identity and subsumption. In R. Green, C. Bean, & S. Myaeng (Eds.), *The semantics of relationships: An interdisciplinary perspective*. Dordrecht: Kluwer.

Guarino, N., & Welty, C. (2002b). Evaluating ontological decisions with ontoclean. *Communications of the ACM, 45*(2), 61–65.

Levin, B. (1993). *English verb classes and alternations: A Preliminary investigation*. Chicago, IL: University of Chicago Press.

Miller, G. A. (1990). WordNet. *Special Issue of the International Journal of Lexicography, 3*(4), 151–161.

Miller, G. A. (1995). WordNet: A lexical database for English. *Communications of the ACM, 38*, 39–40.

Miller, G. A., & Hristea, F. (2006). WordNet nouns: Classes and instances. *Computational Linguistics, 32*(1), 1–3.

Moropa, K., Bosch, S., & Fellbaum, C. (2007). Introducing the African languages WordNet. In *Proceedings of ALASA*, Pretoria, South Africa.

Niles, I., & Pease, A. (2001). Towards a standard upper ontology. In *Proceedings of FOIS-2* (pp. 2–9). Maine: Ogunquit.

Niles, I., & Pease, A. (2003). Linking lexicons and ontologies: Mapping WordNet to the suggested upper merged ontology. In *Proceedings of the international conference on information and knowledge engineering* (pp. 5–6).

Pala, K., Bosch, S., & Fellbaum, C. (2008). Building resources for African languages. In *Proceedings of the sixth international language resources and evaluation*, Marrakech, Morocco.

525  Pease, A., & Fellbaum, C. (2009). Formal ontology as interlingua. In C. R. Huang & L. Prevot (Eds.),
526      *Ontologies and lexical resources.* Cambridge: Cambridge University Press.
527  Robkop, K., Thoongsup, S., Charoenporn, T., Sornlertlamvanich, V., & Isahara, H. (2010). WNMS:
528      Connceting the distributed Wordnet in the case of Asian WordNet. In *The 5th international*
529      *conference of the global WordNet association* (*GWC-2010*), Mumbai, India.
530  Ruppenhofer, J., Baker, C. F., & Fillmore, C. (2002). The FrameNet database and software tools. In
531      A. Braasch & C. Povlsen (Eds.), *Proceedings of the tenth Euralex international congress* (pp. 371–
532      375), Copenhagen, Denmark.
533  Sinha, M., Reddy, M., & Bhattacharyya, P. (2006). An approach towards construction and application of
534      multilingual Indo–WordNet. In *Proceedings of the third global wordnet conference*, Jeju Island,
535      Korea.
536  Tufis, D. (Ed.). (2004). The BalkaNet project. *Special Issue of the Romanian Journal of Information*
537      *Science and Technology, 7*(15), 253–256.
538  Vossen, P. (Ed.). (1998). *EuroWordNet*. Dordrecht: Kluwer.
539  Vossen, P., & Fellbaum, C. (2009). Universals and idiosyncrasies in multilingual wordnets. In H. Boas
540      (Ed.), *Multilingual lexical resources*. Berlin: de Gruyter.
541  Vossen, P., & Rigau, G. (2010). Division of semantic labor in the global wordnet grid. In P. Bhattacharya,
542      C. Fellbaum & P. Vossen (Eds.), *Proceedings of the 5th global WordNet conference*. Narosa
543      Publishing House.
544  Vossen, P., Peters, W., & Gonzalo, J. (1999). Towards a universal index of meaning. In *Proceedings of*
545      *ACL-99 workshop, siglex-99, standardizing lexical resources* (pp. 81–90). University of Maryland,
546      College Park, MD.
547  Vossen, P., Agirre, E., Calzolari, N., Fellbaum, C., Hsieh, S.-K., Huang, C.-R., et al. (2008). Kyoto:
548      A system for mining, structuring, and distributing knowledge across languages and cultures. In
549      *Proceedings of LREC,* Marrakech, Morocco.
550