

Building a WordNet for Arabic

Sabri Elkateb, William Black

The University of Manchester
PO Box 88, Sackville St, Manchester, M60 1QD
w.black@manchester.ac.uk, sabrikom@hotmail.com

Horacio Rodriguez

Politechnical University of Catalonia
Jordi Girona, 1-3, 08034 Barcelona, SPAIN
horacio@lsi.upc.edu

Musa Alkhalifa

University of Barcelona, Edifici Aribau,
5 Planta, Despatx 5.19, Gran via 585,
08007 Barcelona, SPAIN
musa@thera-clic.com

Piek Vossen

Irion Technologies
Irion Technologies, Delftechpark 26, 2628XH,
Delft, The Netherlands
piek.vossen@irion.nl

Adam Pease

Articulate Software Inc, 278 Monroe Dr. #30
Mountain View, CA 94040
apease@articulatesoftware.com

Christiane Fellbaum

Princeton University, Department of Psychology,
Green Hall, Princeton, NJ 08544
fellbaum@clarity.princeton.edu

Abstract

This paper introduces a recently initiated project that focuses on building a lexical resource for Modern Standard Arabic based on the widely used Princeton WordNet for English (Fellbaum, 1998). Our aim is to develop a linguistic resource with a deep formal semantic foundation in order to capture the richness of Arabic as described in Elkateb (2005). Arabic WordNet is being constructed following methods developed for EuroWordNet (Vossen, 1998). In addition to the standard wordnet representation of senses, word meanings are also being defined with a machine understandable semantics in first order logic. The basis for this semantics is the Suggested Upper Merged Ontology and its associated domain ontologies (Niles and Pease, 2001). We will greatly extend the ontology and its set of mappings to provide formal terms and definitions for each synset. Tools to be developed as part of this effort include a lexicographer's interface modeled on that used for EuroWordNet, with added facilities for Arabic script, following Black and Elkateb's earlier work (2004).

Introduction

In recent years, a number of wordnet building efforts have been initiated and carried out within a common framework for lexical representation and are becoming increasingly important resources for a wide range of Natural Language Processing applications. "They can be used in meaning-based information retrieval (searching for concepts rather than specific word forms), in logical inference (if a document mentions dogs, a wordnet allows the inference that it is about animals), in word sense disambiguation (providing the search space of alternative meanings), etc." (Dyvik, 2002). The success of the Princeton WordNet (PWN) for English has motivated similar projects that aim at developing wordnets for other languages. In this paper, we describe our methodology for building a wordnet for Modern Standard Arabic (MSA). This Arabic WordNet (AWN) is to be based on the design and contents of the PWN and can be linked directly to PWN 2.0 and EuroWordNet (EWN), enabling translation

on the lexical level to and from English and dozens of other languages. The Suggested Upper Merged Ontology (SUMO) is being enlarged to provide a formal semantic foundation for AWN (Black et al. 2006). The AWN database will be freely and publicly available.

Challenges

Arabic is a Semitic language which differs from Indo-European languages syntactically, morphologically and semantically. The term 'classical Arabic' refers to the standard form of the language used in all writing and heard on television, radio and in public speeches and religious sermons. The writing system of Arabic has twenty five consonants and three long vowels that are written from right to left and take different shapes according to their position in the word. In addition to the long vowels, Arabic has short vowels. Short vowels are not part of the alphabet but rather are written as vowel diacritics above or under a consonant to give it its desired sound and hence give a word a desired meaning. Texts without vowels are considered to be more appropriate by the Arabic-speaking community since this is the usual form of everyday written and printed materials (books, magazines, newspapers, letters, etc.). But when it comes to the text of the Holy Koran, and more generally to printed collections of classical poetry, school books and some Arabic paper dictionaries, vowel diacritics appear in full. It is very usual for well-edited books, some printed texts, and manuscripts to have vowel diacritics partially or randomly written out in cases where words will be ambiguous or difficult to read. For instance, a word in Arabic consisting of two letters like (ب), i.e., 'b' and 'r', can be very ambiguous without vowel diacritics. Consider the examples in Table 1. Especially in such cases as these, a writer may use diacritics so readers can easily resolve any ambiguity. However, although most Arabs can read texts with vowels explicitly indicated, fewer can write texts using the correct vowel diacritics.

Arabic word	Transliteration	POS	Meaning
بَرّ	barr short vowel 'a'	noun	land (as opposed to sea)
بَرّ	barr short vowel 'a'	adj	reverent, dutiful, kind
بُرّ	burr short vowel 'u'	noun	wheat
بِرّ	birr short vowel 'i'	noun	reverence, kindness

Table 1: vowel diacritics on 'b' and 'r'

For this reason it is a mistake to rely on users, regardless of their background, to correctly enter a search word requiring vowel diacritics. Yet misuse of a single diacritic, such as the 'suku:n' which indicates that a consonant is not followed by any vowel, or as the 'shaddah' (as in *barr* in Table 1 and *darrasa* in Table 2), which indicates a double consonant, will cause a query to fail. People also tend to make mistakes about the position of some diacritics in a word. This can pose a serious problem for information retrieval systems and computerized lexical resources which depend on well-formed user input and may even result in users rejecting the system. In particular, there may be an outright rejection of a robust new lexical resource such as AWN unless that new resource assumes that most of the Arabic speaking users do not have expert command in writing vowel diacritics and will generally ignore them. These users are more comfortable reading texts without diacritics in dealing with everyday written materials including legal and business contracts, newspapers, books as well as both paper and computerized dictionaries. The end result is that it is preferable to allow users to enter Arabic words without diacritics while at the same time allowing the retrieval of those words with vowel diacritics for the purposes of disambiguation.

Another fact about Arabic to take into consideration is that the language has neither capital letters (for proper names: the names of people, countries, cities, geographical features, of months, days of the week, etc.) nor acronyms. This creates increased ambiguity and especially complicates such tasks as Information Extraction in general and Named Entity Recognition in particular.

An additional property of Arabic that should be kept in mind is that Arabic is a highly derivational and inflectional language and its vocabulary can be easily expanded using a framework that is latent in the creative use of roots and morphological patterns. According to Al-Fedaghi and Al-Anzi (1989), cited in De Roeck and Al-Fares (2000), "85% of words derived from tri-literal roots" and there are around 10,000 independent roots.

Because of this, it is possible to build any necessary semantic relation among words of different syntactic categories. That is to say, most Arabic words are created by applying distinct derivational patterns to some root, relating the two not only in form and meaning but determining their syntactic category as well. New Arabic words can always be coined from an existing root

according to the standard derivational patterns. It is also possible to organize sets of Arabic words into distinct semantic fields according to the root from which they are derived. An example of such a field for the root *drs*, 'to study,' is shown in Table 2. Arabic can also adapt loan words from other languages to its system of derivational morphology in order to make them sound and behave like Arabic words as, for example, in the case of *aksadah*, 'oxidation,' which is patterned on *fa'lalah* (Elkateb, 2005).

Arabic word	POS	Pattern	Meaning
دَرَسَ darasa	verb	فَعَّلَ fa?ala	study
دَرَّسَ darrasa	verb	فَعَّلَّ fa??ala	teach
دَرْسٌ dars	noun	فَعْلٌ fa?l	lesson
دِرَاسَةٌ dirasah	noun	فِعْلَالَةٌ fi?alah	study
مُدَرِّسٌ mudarris	noun	مُفَعِّلٌ mufa??il	teacher
مَدْرَسَةٌ madrasah	noun	مَفْعَلَةٌ maf?alah	school
تَدْرِيسٌ tadris	noun	تَفْعِيلٌ taf?il	teaching
تَدَارَسَ tadarasa	verb	تَفَاعَلَ tafa?ala	discuss
دِرَاسِيّ dirasi	adj	فِعْلَالِيّ fi?a:li	educational

Table 2: derivatives of root (d r s)

Numerous efforts have been devoted to the processing of Arabic morphology which outcome is apparent in several approaches and various technical morphological analysers and generators. Among other computational approaches to Arabic morphology, using techniques of Finite State Transducer (FST) and two-level morphology is Beesley (1998, 2001) His system dealt with root, stem and pattern morphology using only two layers. One layer corresponds to the root and is represented by the root lexicon and the other to the morphological measure including vowel pattern.

However, in order to produce a system on the basis of morphological analysis and generation that is linguistically and computationally efficient; the following factors have to be taken into consideration:

1. A word pattern usually combines with a vast number of roots. Roots and patterns are intersected at compile time to yield 90,000 stems. Various combination of prefixes and suffixes, concatenated to the stems, yield over 72,000,000 abstract words.
2. The existence of one morphological form depends on the existence of other forms comprised of the same morphological unit.
3. There are cases where a single form has more than one morphological function as illustrated in Table 1 above.
4. A word is generated by the combination of a root encoded manually and a diacritized pattern each of which has to be hand coded to indicate the subset of patterns with which a root can combine.

5. A root can be extracted by removing the affixes to identify the base form of the diacritized word and to apply it to a morphological measure or a pattern. In this case both word and pattern must be entered manually.
6. Some techniques are designed not to take any Arabic text as an input directly, but to transliterate the Arabic system into ASCII to be fed to the system. The results must be transliterated back to Arabic to be understood. This technique was introduced by Buckwalter (2002) and can be said to have achieved considerable results in Arabic morphological analysis, yet it is unable to adequately deal with ambiguous forms but can only provide full listing of all the possible readings of the ambiguous form.

There seems to be no agreement on the nearest way to adequate morphological analysis/generation and there is yet no proper means for generating or analyzing the Arabic roots due to the complexity of the weak vowels governing a vast amount of the vocabulary. It seems also that there is no role for morphological generation in suggesting words, because for much of the vocabulary, the rate at which these would prove to be actual words would be too low unless at least three quarters of the process are done manually (Elkateb, 2005). As far as dictionaries are concerned, a multilingual resource generally includes equivalence and translation relations and should tackle issues like language specific and untranslatable material. Translation is not merely an act of linguistic transfer, but it also involves the interaction of cultures and that transference of culture imposes far greater problems than linguistic transfer. Translation of words of cultural content may involve solving problems like the unavailability of equivalents or tackling untranslatable items and consequently filling the gaps that may exist among languages. Consider the Arabic words in Table 3

<i>zaka:t</i>	annual compulsory alms (2.5 %) of the savings of a Muslim when any amount or property exceeds one year in possession.
<i>suhu:r</i>	a light meal before starting a new fasting day of Ramadan (before daybreak).
<i>hija:b</i>	an Islamic veil which is worn by women to cover the hair and the neck.
mu'akhar Sada:q	money/property stipulated upon in the marriage contract which is due to be paid by the husband to his wife in case he intends to divorce her.

Table 3: lexical gaps

Lexical Ambiguity

A lexical item may carry two distinct and unrelated meanings, i.e. homonymy. A homonym can be defined as a word with no relationship between its senses, as in the word *bank* where the first sense refers to a river side and the second to a financial institution. Ambiguity and polysemy of nominal forms represent an important concern which affects the organization of word meaning.

The basic distinction between what Pustejovsky, (1995) termed contrastive ambiguity and complementary polysemy should involve different solutions for the representation of lexical knowledge. Contrastive ambiguity, as manifested by words such as *bank* (financial institution or river side) is handled by multiple representations for the clarity of senses. However it is claimed that this type does not form a significant problem in the language since contrastive ambiguity between two unrelated senses of a word tends to be a historically accidental and idiosyncratic property of individual words. Hence “we don’t expect to find instances of the same contrastive ambiguity replicated by other words in the language or by words in other languages” (Dyvik, 2003). Complementary polysemy occurs in cases where a single word has multiple senses which are related to one another in some predictable way. It is claimed that ambiguity can result from senses which are manifestations of the same basic meaning of the word depending on the context it occurs in. The manner in which senses are related in complementary polysemy is the factor that distinguishes it from contrastive ambiguity where senses have no contextual relation. Accordingly, a word like ‘*door*’ has two related senses being (physical object or aperture). So, knocking on the ‘*door*’ (physical object) is different from going through the same ‘*door*’ (aperture). Let us first examine the senses of the Arabic word ‘*bab*’ for ‘*door*’ in order to figure out how words behave in different languages and how sense extensions vary from one language to another:

bab (door/chapter)

--sense₁ = *physical object*, e.g. I painted the front door.

--sense₂ = *aperture* e.g. Adam went through the door.

--sense₃ = *written communication (book chapter)*, “opening of a piece of text” e.g. I started a new chapter of my thesis.

The first two senses are more closely related than the third. The third sense in Arabic refers to opening/entering (or going through writing/reading) a written text. This sense might be extended from the notion of ‘*opening*’ as in ‘*open the book*’ or ‘*open a new chapter*’ compared to ‘*open the door*’. Therefore, it can be said to be an instance of complementary polysemy not contrastive ambiguity because of the shared collocates with the verb to open.

It is claimed that complementary polysemy poses a serious problem not only in one language but also would normally be projected into other languages. The English word ‘*lamb*’, for example, is said to denote two different senses: a count noun *animal* and a mass noun *meat* whereas in Arabic the word ‘*hamal*’ (*lamb*) and its synonyms ‘*kharu:f*’ (*lamb/sheep*) refer only to the count noun ‘*animal*’. It seems that it is only accidentally, in English, that this noun is classified as polysynonymous because it refers to both *animal* and *meat*. This may be because it is linked with small masses like ‘*chicken, eggs, snails*’ where complementary polysemy is less frequent. More

interestingly, the polysemy in the case of *lamb* is only temporary and will disappear as the lamb gets old and becomes a sheep. The second sense for ‘*lamb*’ as mass noun ‘*meat*’ can only appear in Arabic if the word *lamb* occurs in a compound as in ‘*lahm kharu:f*’ (*sheep meat/mutton*) where the complementary polysemy is completely absent. However, Arabic and English interpret other masses the same way whether large or small, like ‘*fish*’, ‘*chicken*’, ‘*eggs*’, ‘*potatoes*’ etc., where complementary polysemy may occur equally in both languages:

1. I did not like the fish we had for lunch.
2. I went to see the dead fish at lunch time.

There are cases in Arabic where a word may carry multiple but related senses as in the noun ‘*sawt/aswat*’ where it can be classified as complementary polysemy according to its interpretation in Arabic:

sawt / aswat

--sense₁ = vote: an *indication* of a choice or opinion that is made by voting

--sense₂ = voice: sound produced by *speaking* or singing.

The common morphological derivation of a pair of nouns in Arabic provides evidence for their relatedness as polysemes. The Arabic word ‘*sawt*’ (*vote*) and ‘*swat*’ (*voice*) are apparently derived from the same unaugmented trilateral root ‘*s w t*’ (*sound*). In addition, the ‘*indication*’ of vote in sense₁ refers to verbal consent ‘*speaking*’ in sense₂.

3. *hada fariq ?add al aswat* (This is a vote counting team).
4. *hada fariq tasji:l al aswat* (This is a voice recording team).

The two senses in 3 and 4 can be classified as complementary polysemy rather than contrastive senses i.e., to ‘*vote*’ is to primarily ‘*say*’ who or what you are in favour of. Example 4 above also shows that the word ‘*aswat*’ denotes two senses: ‘*votes*’ and ‘*voices*’ as unrelated to one another when modified by ‘*tasji:l*’ (*recording*) which denotes the recording of voice as well as writing down (in a record) the names of the voters (*votes*). Therefore example 4 can be interpreted as having these two contrastive senses in 5:

5. *hada fariq tasji:l al aswat*:
 - a. This is a voice recording team. (audio recording)
 - b. This is a vote recording team. (writing)

This word gets even more ambiguous in its proper context than on its own or in a lexicon as in example 6:

6. *hadihi aswat alnakhibi:n*.

The word ‘*aswat*’ in this context refers to two different senses:

- a. These are the voices of the electors.
- b. These are the votes of the electors.

Ambiguity varies between two languages when one borrows a word from the other. In this case, polysemy projects into the borrowing language from the source language but not the opposite. The term ‘*alqaida*’ borrowed from Arabic to refer to a group of extremists in Afghanistan known by this name and classified as a terrorist organization. This proper name of this entity is derived from the meaning of ‘*the base*’. Since proper names are not translated, as illustrated in example 7 below, the polysemy in this case occurs only in Arabic but not in English. In other words, the sentence ‘*The Americans attacked Alqaida*’ carries one sense in English whereas in Arabic is interpreted as having two senses:

7. *alamrica:n yuha:jimu:n alqaida*.
 - a. The Americans attacked Alqaida. (terrorist group based in Afghanistan)
 - b. The Americans attacked the base. (a military base)

No one would argue about the importance of a semantic lexicon to handle such different and/or related senses of words and concepts. However, there should be an agreement on how to represent lexical data to be easily manipulated by computers in order to encode any semantic relations between senses and to carry out various applications of a conceptual lexicon such as word sense disambiguation (WSD), lexical chains etc.

Lexicography

Following EuroWordNet, AWN is developed in two phases by first building a core wordnet around the most important concepts, the so-called Base Concepts (Vossen 1998), and secondly extending the core wordnet downward to more specific concepts using additional criteria. The core wordnet should thus become highly compatible with wordnets in other languages that are developed according to the same approach.

For the core wordnet, The Common Base Concepts (CBCs) of the 12 languages in EWN and BalkaNet (Tufis, 2004) are being encoded as synsets in AWN; other Arabic language-specific concepts are added and translated manually to the closest synset. The same procedure is performed for all English synsets that currently have an equivalence relation in the SUMO ontology. Synset encoding proceeds bi-directionally: given an English synset, all corresponding Arabic variants (if any) will be selected; given an Arabic word, all its senses are determined and for each of them the corresponding English synset is encoded.

The Arabic synsets will be extended with hypernym relations to form a closed semantic hierarchy. SUMO will be used to maximize the semantic consistency of the hyponymy links. This will represent the core wordnet, which is a semantic basic for the further extension. The work is mostly done manually.

When a new Arabic verb is added, extensions are made from verbal entries, including verbal derivatives,

nominalizations, verbal nouns, and so on. We also consider the most productive forms of deriving broken plurals. This is done by applying lexical and morphological rules iteratively.

The database is further extended downward from the CBCs. First, a layer of hyponyms is chosen based on maximal connectivity, relevance, and generality. Two major pre-processing steps are required, preparation and extension. Preparation entails compiling lexical and morphological rules and processing available bilingual resources from which we construct a homogeneous bilingual dictionary containing information on the Arabic/English word pair. This information includes the Arabic root, the POS, the relative frequencies and the sources supporting the pairing. The Arabic words in these bilingual resources must also be normalized and lemmatized while maintaining vowels and diacritics.

We next apply 17 heuristic procedures, previously used for EWN, to the bilingual dictionary in order to derive candidate Arabic words/English synsets mappings. Each mapping includes the Arabic word and root, the English synset, the POS, the relative frequencies, a mapping score, the absolute depth in AWN, the number of gaps between the synset and the top of the AWN hierarchy, and attested tokens of the pair. The Arabic word/English synset pairs constitute the input to a manual validation process. We proceed by chunks of related units (sets of related WN synsets, e.g. hyponymy chains and sets of related Arabic words, i.e., words having the same root) instead of individual units (i.e., synsets, senses, words).

Finally, AWN will be completed by filling in the gaps in its structure, covering specific domains, adding terminology and named entities, etc. Each synset construction step is followed by a validation phase, where formal consistency is checked and the coverage is evaluated in terms of frequency of occurrence and domain distribution. The total coverage of AWN will be around 10,000 synsets.

Tools

Tools to be developed for AWN include a lexicographer's interface modeled on the EWN interface with added facilities for Arabic script. Because AWN is to be aligned not just to PWN but to every wordnet aligned to PWN – either directly or indirectly through an Interlingual Index or the ontology – the database design supports multiple languages. The user interface will be explicitly multilingual and indifferent to the direction of alignment between the conceptual structures of the two languages. In addition to search and browsing facilities for the end users of the completed database, lexicographers require an editing interface. A variety of legacy components are available, each with their relative advantages. The editor's interface will communicate with the database server using Simple Object Access Protocol (SOAP), allowing multiple lexicographers at different sites to maintain a common database.

Database

The database structure comprises four principal entity types: *item*, *word*, *form* and *link*. *Items* are conceptual

entities, including synsets, ontology classes and instances. An item has a unique identifier and descriptive information such as a gloss. Items lexicalized in different languages are distinct. A *word* entity is a word sense, where the word's citation form is associated with an item via its identifier. A *form* is an entity that contains lexical information (not merely inflectional variation). The forms are the root and/or the broken plural form, where applicable. A *link* relates two items, and has a *type* such as "equivalence," "subsuming," etc. Links interconnect sense items, e.g., a PWN synset to an AWN synset, a synset to a SUMO concept, etc. This data model has been specified in XML as an interchange format, but is also implemented in a MySQL database hosted by one of the partners.

Ontology

A large ontology providing the semantic underpinning for AWN concepts will be built on SUMO, a formal ontology of about 1000 terms and 4000 definitional statements currently that is provided in a first order logic language called Standard Upper Ontology Knowledge Interchange format (SUO-KIF) and also translated into OWL semantic web language. SUMO has natural language generation templates and a multi-lingual lexicon that allows statements in SUO-KIF and SUMO to be expressed in multiple languages. Synsets map to a general SUMO term or a term that is directly equivalent to the given synset (Figure 1).

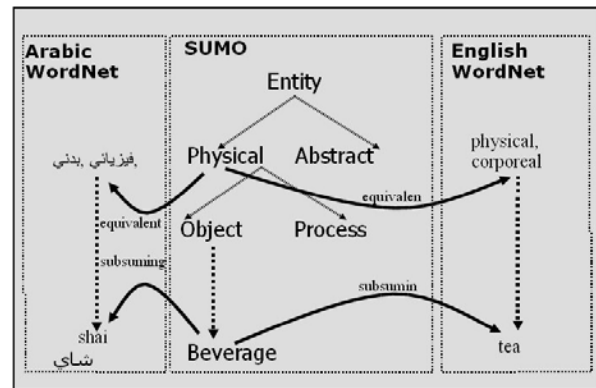


Figure 1: SUMO mapping to wordnets

New formal terms will be defined to cover a greater number of equivalence mappings, and the definitions of the new terms will in turn depend upon existing fundamental concepts in SUMO. The process of formalizing definitions will generate feedback as to whether word senses in AWN need to be divided or combined and how glosses may be clarified. Wordnets in other languages linked by synset number will benefit, too. The Sigma ontology development environment will be updated to handle a similar presentation of Unicode-based character sets, including Arabic.

The Interlingual Index (ILI) connecting EWN wordnets is a condensed set of more or less universal concepts linking synsets across languages via multiple exhaustive equivalence relations. In EuroWordNet and BalkaNet,

English PWN has been used to express equivalence relations across the different languages. By providing many SUMO definitions and terms that correspond to Arabic synsets, we will create the opportunity to use SUMO as the ILI for all wordnets that are currently related to PWN. This is illustrated in Figure 2. If the Arabic word sense for *shai* is exhaustively defined by relations to SUMO terms, this definition can replace an equivalence relation (er1) that is currently encoded between the Arabic synset *shai* and a synset *tea* in PWN. Note that the relations from *shai* to the SUMO terms need to be exhaustive, which may require multiple relations of different types (sr1 (subsumption), r2, r3) to multiple SUMO terms.

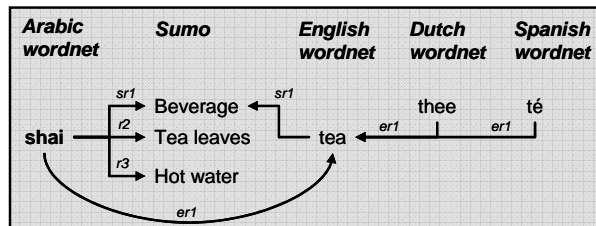


Figure 2: SUMO and ILI

If there are also equivalence relations from other languages (e.g. Dutch and Spanish) to the same PWN synset, then these relations grant the linkage of the synsets in these languages to the same SUMO definition.

Besides providing a formal semantic framework, SUMO can thus also be used to map synsets across languages, in fact even when there is not an equivalent in English. By composing formal definitions for the non-English synsets, SUMO as an ILI will not only be less biased by English but also has more expressive power.

Conclusion

Constructing AWN presents challenges not encountered by established wordnets. These include the script on the one hand and the morphological properties of Semitic languages, centered around roots, on the other hand. The foundations for meeting these challenges have been laid. An innovation with significant consequences for wordnet development is the proposal to substitute English WN as the ILI with SUMO.

Acknowledgements

This work was supported by the United States Central Intelligence Agency.

References

Beesley, K. (2001) Finite-State Morphological Analysis and Generation of Arabic at Xerox, ACL/EACL 2001, July 6th, Toulouse, France : 1-8
 Black, W., Elkateb, S., Rodriguez, H., Alkhalifa, M., Vossen, P., Pease, A. and Fellbaum, C., (2006). Introducing the Arabic WordNet Project, in Proceedings

of the Third International WordNet Conference, Sojka, Choi, Fellbaum and Vossen eds.
 Black, W. J., and Elkateb, S. (2004) A Prototype English-Arabic Dictionary Based on WordNet, Proceedings of 2nd Global WordNet Conference, GWC2004, Czech Republic, 67-74.
 Buckwalter, T. (2002) Arabic Morphological Analysis, [Http://www.qamus.org/morphology.htm](http://www.qamus.org/morphology.htm)
 De Roeck, A., and Al-Fares, W. (2000) A Morphologically Sensitive Clustering Algorithm for Identifying Arabic Roots Proceedings of the 38th Annual Meeting of the ACL, Hong Kong, 199-206
 Dyvik, H. (2003) Translations as a semantic knowledge source: word alignment and wordnet, Section for Linguistic Studies scientific papers, University of Bergen
 Dyvik, H. (2002) Translations as Semantic Mirrors: From Parallel Corpus to Wordnet1. Section for Linguistic Studies scientific papers, University of Bergen
 Elkateb, S and Black, W. J. (2001) Towards the Design of English-Arabic Terminological Knowledge Base, Proceedings of ACL 2000, Toulouse, France:113-118
 Elkateb, S and Black, W. J. (2004) A Bilingual Dictionary with Enriched Lexical Information, Proceedings of NEMLAR Cairo, Egypt 2004 Arabic Language Tools and Resources: 79-84
 Elkateb, S. (2005) Design and implementation of an English Arabic dictionary/editor. PhD thesis, The University of Manchester, United Kingdom.
 Farreres, J. (2005) Creation of wide-coverage domain-independent ontologies. PhD thesis, Universitat Politècnica de Catalunya.
 Fellbaum, C., (1998, ed.) WordNet: An Electronic Lexical Database. Cambridge, MA: MIT Press.
 Niles, I., and Pease, A. (2001) Towards a Standard Upper Ontology. In: Proceedings of FOIS 2001, Ogunquit, Maine, pp. 2-9.
 Pease, A., (2000) Standard Upper Ontology Knowledge Interchange Format. Web document <http://suo.ieee.org/suo-kif.html>.
 Pease, A., (2003) The Sigma Ontology Development Environment, in Working Notes of the IJCAI-2003 Workshop on Ontology and Distributed Systems. Volume 71 of CEUR Workshop Proceeding series
 Pustejovsky, J. (1995) The Generative Lexicon, Massachusetts Institute of Technology.
 Tufis, D. (ed.) (2004) Special Issue on the BalkaNet project. Romanian Journal of Information Science and Technology, Vol.7, nos 1-2
 Vossen, P. (ed.) (1998) EuroWordNet: A Multilingual Database with Lexical Semantic Networks. Dordrecht: Kluwer Academic Publishers.
 Vossen P. (2004) EuroWordNet: a multilingual database of autonomous and language-specific wordnets connected via an Inter-Lingual-Index. International Journal of Lexicography, Vol.17 No. 2, OUP, 161-173