



## Consistent Criteria for Sense Distinctions

MARTHA PALMER

*Department of Computer Science, IRCS, University of Pennsylvania, Phil, PA 19104, USA*  
(E-mail: [mpalmer@cis.upenn.edu](mailto:mpalmer@cis.upenn.edu))

**Abstract.** This paper specifically addresses the question of polysemy with respect to verbs, and whether or not the sense distinctions that are made in on-line lexical resources such as WordNet are appropriate for computational lexicons. The use of sets of related syntactic frames and verb classes are examined as a means of simplifying the task of defining different senses, and the importance of concrete criteria such as different predicate argument structures, semantic class constraints and lexical co-occurrences is emphasized.

### 1. Introduction

The difficulty of achieving adequate hand-crafted semantic representations has limited the field of natural language processing to applications that can be contained within well-defined subdomains. The only escape from this limitation will be through the use of automated or semi-automated methods of lexical acquisition. However, the field has yet to develop a clear consensus on guidelines for a computational lexicon that could provide a springboard for such methods, in spite of all of the effort on different lexicon development approaches (Mel'cuk, 1988; Pustejovsky, 1991; Nirenburg et al., 1992; Copestake and Sanfilippo, 1993; Lowe et al., 1997; Dorr, 1997). One of the most controversial areas has to do with polysemy. What constitutes a clear separation into senses for any one verb or noun, and how can these senses be computationally characterized and distinguished? The answer to this question is the key to breaking the bottleneck of semantic representation that is currently the single greatest limitation on the general application of natural language processing techniques.

In this paper we specifically address the question of polysemy with respect to verbs, and whether or not the sense distinctions that are made in on-line dictionary resources such as WordNet (Miller, 1990; Miller and Fellbaum, 1991), are appropriate for computational lexicons. We examine the use of sets of related syntactic frames and verb classes as a means of simplifying the task of defining different senses, and we focus on the mismatches between these types of distinctions and some of the distinctions that occur in WordNet.

## 2. Challenges in Building Large-Scale Lexicons

Computational lexicons are an integral part of any natural language processing system, and perform many essential tasks. Machine Translation (MT), and Information Retrieval (IR), both rely to a large degree on isolating the relevant senses of words in a particular phrase, and there is wide-spread interest in whether or not word sense disambiguation (WSD), can be performed as a separate self-contained task that would assist these applications.<sup>1</sup> Information retrieval mismatches such as the retrieval of an article on plea bargaining, (*speedier trials and lighter sentences*), given *speed of light* as a query are caused by inadequate word sense disambiguation. These are clearly not the same senses of *light*, (or even the same parts of speech), but a system would have to distinguish between WordNet light1, involving visible light, and WordNet light2, having to do with quantity or degree in order to rule out this retrieval. However, it is possible that the lexically based statistical techniques currently employed in the best IR systems are already accomplishing a major portion of the WSD task, and a separate WSD stage would have little to add (Voorhees, 1999). Clear sense distinctions have a more obvious payoff in MT. For instance, in Korean, there are two different translations for the English verb *lose*, depending on whether it is an object that has been misplaced or a competition that has been lost: lose1, *lose the report – ilepeli-ess-ta*, and lose2, *lose the battle – ci-ess-ta* (Palmer et al., 1998). Whether or not WSD is a useful separate stage of processing for MT or part of an integrated approach, selecting the appropriate entry in a bilingual lexicon is critical to the success of the translation.

The *lose* sense distinctions can be made by placing semantic class constraints on the object positions, i.e., +competition, and +solid object respectively. The first constraint corresponds directly to a WordNet hypernym, but the second one does not. The closest correlate in WordNet would be +abstract activity, which is the common hypernym for both *hostile military engagement* and *game*, and which may discriminate sufficiently.

Computational lexicons can most readily make sense distinctions based on concrete criteria such as:

- different predicate argument structures
- different semantic class constraints on verb arguments
- different lexical co-occurrences, such as prepositions

This seems straightforward enough, and traditional dictionaries usually have separate entries for transitive (two argument) and intransitive (one argument) verbs, as well as for verb particle constructions (with specific prepositions, as in *break off*). However, semantic class constraints are never made explicit in dictionaries, and lexicographers often refer to even more inaccessible implicit criteria. For instance, out of the ten senses that WordNet 1.6, gives for *lose*, we find one, WN2, that corresponds to our lose1 from above, *lose the battle* sense, but two, WN1 and WN5, that correspond to our lose2, *misplace an item*.

- lose1 – WN2. lose – (fail to win; “We lost the battle but we won the war”)
- lose2 – WN1. (fail to keep or to maintain; cease to have, either physically or in an abstract sense; fail to keep in mind or in sight; “She lost her purse when she left it unattended on her seat”; “She lost her husband a year ago”)
- lose2 – WN5. (miss from one’s possessions; lose sight of; “I’ve lost my glasses again!”)

When we try to establish concrete criteria for distinguishing between WN1 (*lost her purse*) and WN5 (*lost my glasses*), we realize that these two WordNet senses are not distinguished because of anything to do with semantic class constraints on the verb arguments (an +animate Agent and a +solid object possessed by the Agent in both cases), but rather are distinguished by possible future events – namely the likelihood of the object being found. It is not reasonable to expect a computational lexicon to characterize all possible worlds in which an event can take place, and then distinguish between all possible outcomes. A more practical sense division for a computational lexicon would be [lose1 (losing competitions), lose2 (misplacing objects), lose3 (being bereft of loved ones)].<sup>2</sup>

We are not denying that a computational lexicon should include particular changes in the state of the world that are entailed by specific actions, quite the contrary (Palmer, 1990). However, the characterizations of these changes should be generally applicable, and cannot be so dependent on a single world context that they change with every new situation.

Other areas of difference between computational lexicons and more traditional lexical resources have to do with the flexibility of the representation. Computational lexicons are particularly well suited to capturing hierarchical relationships and regular sense extensions based on verb class membership. For instance, the following two senses are in and among the 63 sense distinctions WordNet listed for *break*.

- break – WN2. break, separate, split up, fall apart, come apart – (become separated into pieces or fragments; “The figurine broke”; “The freshly baked loaf fell apart”)
- break – WN5. (destroy the integrity of; usually by force; cause to separate into pieces or fragments; “He broke the glass plate”; “She broke the match”)

They are shown as being related senses in WordNet 1.6, but the relationship is not made explicit. It is a simple task for a computational lexicon to specify the type of relationship, i.e., the transitive frame in WN5 the causative form of WN2, and has explicit inclusion of an Agent as an additional argument. In the XTAG English lexicon (Joshi et al., 1975; Vijay-Shanker, 1987), this is currently handled by associating both the intransitive/ergative and transitive tree families<sup>3</sup> with the same syntactic database entry for *break*. In the transitive form the NP1 (Patient) becomes the Object and an NP0 (Agent) is added as the Subject. The +**causative**

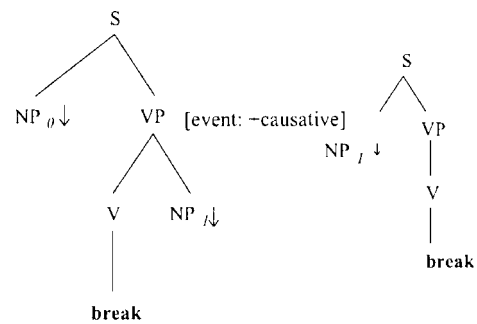


Figure 1. An ergative verb and its causative sense extension.

semantic feature can be added as well, as illustrated in Figure 1.<sup>4</sup> We are currently adding syntactic frames to the two related entries in WordNet 1.6, which, as well as making the definitions more consistent, helps to explicitly capture the sense extension. This resource, called VerbNet, will be available soon (Dang et al., 1998).

In addition to regular extensions in meaning that derive from systematic changes in subcategorization frames, there are also regular extensions occasioned by the adjunction of optional prepositions, adverbials and prepositional phrases. For example, the basic meaning of *push*, *He pushed the next boy*, can be extended to explicitly indicate accompanied motion by the adjunction of a path prepositional phrase, as in *He pushed the boxes across the room* (Palmer et al., 1997; Dang et al., 1998), which corresponds to WN1 below. The possibility of motion of the object can be explicitly denied through the use of the conative, as in *He pushed at the box*, which is captured by WN5. Finally, the basic sense can also be extended to indicate a change of state of the object by the adjunction of *apart*, as in *He pushed the boxes apart*. There is no WordNet sense that corresponds to this, nor should there be. What is important is for the lexicon to provide the capability of recognizing and generating these usages where appropriate. If they are general enough to apply to entire classes of verbs, then they can be captured through regular adjunctions rather than being listed explicitly (for more details, see Bleam et al., 1998).

- WN1. push, force – (move with force, “He pushed the table into a corner”; “She pushed her chin out”)
- WN5. push – (press against forcefully without being able to move)

### 3. Conclusion

It has been suggested that WordNet sense distinctions are too fine-grained and coarser senses are needed to drive the word sense disambiguation task. For instance, in defining *cut*, WordNet distinguishes between WN1, *separating into pieces of a concrete object*, WN29, *cutting grain*, WN30, *cutting trees*, and WN33,

*cutting hair*. For many purposes, the three more specialized senses, WN29, WN30 and WN33, which all involve *separation into pieces of concrete objects* could be collapsed into the more coarse-grained WN1. However, when searching for articles on recent changes in hair styles, the more fine-grained WN33 would still be useful. Computational lexicons actually lend themselves readily to moving back and forth between elements of an hierarchical representation based on concrete criteria, and this type of structuring should become more prevalent. The point is that they operate most effectively in the realm of concrete criteria for sense distinctions, such as changes in argument structure, changes in sets of syntactic frames and/or changes in semantic class constraints, and lexical co-occurrences. Distinctions that are based on world knowledge, no matter how diverse, are much more problematic. We must bear this in mind in order to design a word sense disambiguation task that will also encourage rational, incremental development of computational lexicons.

### Acknowledgements

We thank Aravind Joshi, the members of the XTAG group, Christiane Fellbaum and our reviewers. This work has been supported in part by NSF grants SBR 8920230 and IIS-9800658 and Darpa grant #N66001-94C-6043.

### Notes

<sup>1</sup> For a discussion of WSD and IR, see Krovetz and Croft (1992) and Sanderson (1994).

<sup>2</sup> Obviously, semantic class constraints on the *object* would fail to distinguish between losing one's husband in the supermarket versus losing one's spouse to cancer, and additional information such as adjuncts would have to be considered as well.

<sup>3</sup> A tree family contains all of the syntactic realizations associated with a particular subcategorization frame, such as subject and object extraction and passive (XTAG-Group, 1995; Xia et al., 1999).

<sup>4</sup> All of Levin's *break* and *bend* verbs are given the same type of entry, as well as many other verbs (Levin, 1993; Dang et al., 1998).

### References

- Bleam, T., M. Palmer and V. Shanker. "Motion Verbs and Semantic Features in Tag". In *Proceedings of the TAG+98 Workshop*. Philadelphia, PA, 1998.
- Copestake, A. and A. Sanfilippo. "Multilingual Lexical Representation". In *Proceedings of the AAAI Spring Symposium: Building Lexicons for Machine Translation*. Stanford, California, 1993.
- Dang, H.T., K. Kipper, M. Palmer and J. Rosenzweig. "Investigating Regular Sense Extensions Based on Intersective Levin Classes". In *Proceedings of Coling-ACL98*. Montreal, CA, 1998.
- Dorr, B.J. "Large-Scale Dictionary Construction for Foreign Language Tutoring and Interlingual Machine Translation". *Machine Translation*, 12 (1997), 1–55.
- Joshi, A.K., L. Levy and M. Takahashi. "Tree Adjunct Grammars". *Journal of Computer and System Sciences* (1975).
- Krovetz, R. and W. Croft. "Lexical Ambiguity and Information Retrieval". *ACM Transactions on Information Systems*, 10(2) (1992), 115–141.
- Levin, B. *English Verb Classes and Alternations: A Preliminary Investigation*. Chicago, IL: The University of Chicago Press, 1993.

- Lowe, J., C. Baker and C. Fillmore. "A Frame-Semantic Approach to Semantic Annotation". In *Proceedings 1997 Siglex Workshop/ANLP97*. Washington, D.C., 1997.
- Mel'cuk, I.A. "Semantic Description of Lexical Units in an Explanatory Combinatorial Dictionary: Basic Principles and Heuristic Criteria". *International Journal of Lexicography*, 1(3) (1988), 165–188.
- Miller, G.A. "Wordnet: An On-Line Lexical Database". *International Journal of Lexicography*, 3 (1990), 235–312.
- Miller, G.A. and C. Fellbaum (1991). "Semantic Networks of English". *Lexical and Conceptual Semantics, Cognition Special Issue*. 1991, pp. 197–229.
- Nirenburg, S., J. Carbonell, M. Tomita and K. Goodman *Machine Translation: A Knowledge-Based Approach*. San Mateo, California, USA: Morgan Kaufmann, 1992.
- Palmer, M. "Customizing Verb Definitions for Specific Semantic Domains". *Machine Translation*, 5 (1990).
- Palmer, M., C. Han, F. Xia, D. Egedi and J. Rosenzweig. "Constraining Lexical Selection Across Languages Using Tags". In *Tree Adjoining Grammars*. Ed. A. Abeille and O. Rambow. Palo Alto, CA: CSLI, 1998.
- Palmer, M., J. Rosenzweig, H. Dang and F. Xia. "Capturing Syntactic/Semantic Generalizations in a Lexicalized Grammar". *Presentation in Working Session of Semantic Tagging Workshop, ANLP-97*. 1997.
- Pustejovsky, J. "The Generative Lexicon". *Computational Linguistics*, 17(4) (1991).
- Sanderson, M. "Word Sense Disambiguation and Information Retrieval". In *Proceedings of the 17th ACM SIGIR Conference*. 1994, pp. 142–151.
- Vijay-Shanker, K. (1987). *A Study of Tree Adjoining Grammars*. PhD thesis, Department of Computer and Information Science, University of Pennsylvania.
- Voorhees, E.M. "Natural Language Processing and Information Retrieval". In *Proceedings of Second Summer School on Information Extraction*, Lecture Notes in Artificial Intelligence. Springer-Verlag, 1999.
- Xia, F., M. Palmer and K. Vijay-Shanker. "Towards Semi-Automating Grammar Development". In *Proceedings of the Natural Language Processing Pacific Rim Symposium (NLPRS-99)*. Beijing, China, 1999.
- XTAG-Group, *A Lexicalized Tree Adjoining Grammar for English*, Technical Report IRCS 95-03. University of Pennsylvania, 1995.