

Mapping Lexical Entries in a Verbs Database to WordNet Senses

Rebecca Green[†] and Lisa Pearl[†] and Bonnie J. Dorr^{†§} and Philip Resnik^{†§}

§Institute for Advanced Computer Studies

†Department of Computer Science

University of Maryland

College Park, MD 20742 USA

{rgreen, llsp, bonnie, resnik}@umiacs.umd.edu

Abstract

This paper describes automatic techniques for mapping 9611 entries in a database of English verbs to WordNet senses. The verbs were initially grouped into 491 classes based on syntactic features. Mapping these verbs into WordNet senses provides a resource that supports disambiguation in multilingual applications such as machine translation and cross-language information retrieval. Our techniques make use of (1) a training set of 1791 disambiguated entries, representing 1442 verb entries from 167 classes; (2) word sense probabilities, from frequency counts in a tagged corpus; (3) semantic similarity of WordNet senses for verbs within the same class; (4) probabilistic correlations between WordNet data and attributes of the verb classes. The best results achieved 72% precision and 58% recall, versus a lower bound of 62% precision and 38% recall for assigning the most frequently occurring WordNet sense, and an upper bound of 87% precision and 75% recall for human judgment.

1 Introduction

Our goal is to map entries in a lexical database of 4076 English verbs automatically to WordNet senses (Miller and Fellbaum, 1991), (Fellbaum, 1998) to support such applications as ma-

chine translation and cross-language information retrieval. For example, the verb *drop* is multiply ambiguous, with many potential translations in Spanish: *bajar*, *caerse*, *dejar caer*, *derribar*, *disminuir*, *echar*, *hundir*, *soltar*, etc. The database specifies a set of interpretations for *drop*, depending on its context in the source-language (SL). Inclusion of WordNet senses in the database enables the selection of an appropriate verb in the target language (TL). Final selection is based on a frequency count of WordNet senses across all classes to which the verb belongs—e.g., *disminuir* is selected when the WordNet sense corresponds to the meaning of *drop* in *Prices dropped*.

Our task differs from standard word sense disambiguation (WSD) in several ways. First, the words to be disambiguated are entries in a lexical database, not tokens in a text corpus. Second, we take an “all-words” rather than a “lexical-sample” approach (Kilgarriff and Rosenzweig, 2000): All words in the lexical database “text” are disambiguated, not just a small number for which detailed knowledge is available. Third, we replace the contextual data typically used for WSD with information about verb senses encoded in terms of thematic grids and lexical-semantic representations from (Olsen et al., 1997). Fourth, whereas a single word sense for each token in a text corpus is often assumed, the absence of sentential context leads to a situation where several WordNet senses may be equally appropriate for a database entry. Indeed, as distinctions between WordNet senses can be fine-grained (Palmer, 2000), it may be unclear, even in context, which sense is meant.

The verb database contains mostly syntactic in-

formation about its entries, much of which applies at the class level within the database. WordNet, on the other hand, is a significant source for information about semantic relationships, much of which applies at the “synset” level (“synsets” are WordNet’s groupings of synonymous word senses). Mapping entries in the database to their corresponding WordNet senses greatly extends the semantic potential of the database.

2 Lexical Resources

We use an existing classification of 4076 English verbs, based initially on *English Verbs Classes and Alternations* (Levin, 1993) and extended through the splitting of some classes into subclasses and the addition of new classes. The resulting 491 classes (e.g., “Roll Verbs, Group I”, which includes *drift*, *drop*, *glide*, *roll*, *swing*) are referred to here as *Levin+ classes*. As verbs may be assigned to multiple Levin+ classes, the actual number of entries in the database is larger, 9611.

Following the model of (Dorr and Olsen, 1997), each Levin+ class is associated with a *thematic grid* (henceforth abbreviated θ -grid), which summarizes a verb’s syntactic behavior by specifying its predicate argument structure. For example, the Levin+ class “Roll Verbs, Group I” is associated with the θ -grid [th goal], in which a theme and a goal are used (e.g., *The ball dropped to the ground*).¹ Each θ -grid specification corresponds to a *Grid class*. There are 48 Grid classes, with a one-to-many relationship between Grid and Levin+ classes.

WordNet, the lexical resource to which we are mapping entries from the lexical database, groups synonymous word senses into “synsets” and structures the synsets into part-of-speech hierarchies. Our mapping operation uses several other data elements pertaining to WordNet: semantic relationships between synsets, frequency data, and syntactic information.

Seven semantic relationship types exist between synsets, including, for example, antonymy, hyperonymy, and entailment. Synsets are often related to a half dozen or more other synsets; they

¹There is also a Levin+ class “Roll Verbs, Group II” which is associated with the θ -grid [th particle(down)], in which a theme and a particle ‘down’ are used (e.g., *The ball dropped down*).

may be related to multiple synsets through a single relationship or may be related to a single synset through multiple relationship types.

Our frequency data for WordNet senses is derived from SEMCOR—a semantic concordance incorporating tagging of the Brown corpus with WordNet senses.²

Syntactic patterns (“frames”) are associated with each synset, e.g., Somebody ___s something; Something ___s; Somebody ___s somebody into V-ing something. There are 35 such verb frames in WordNet and a synset may have only one or as many as a half dozen or so frames assigned to it.

Our mapping of verbs in Levin+ classes to WordNet senses relies in part on the relation between thematic roles in Levin+ and verb frames in WordNet. Both reflect how many and what kinds of arguments a verb may take. However, constructing a direct mapping between θ -grids and WordNet frames is not possible, as the underlying classifications differ in significant ways. The correlations between the two sets of data are better viewed probabilistically.

Table 1 illustrates the relation between Levin+ classes and WordNet for the verb *drop*. In our multilingual applications (e.g., lexical selection in machine translation), the Grid information provides a context-based means of associating a verb with a Levin+ class according to its usage in the SL sentence. The WordNet sense possibilities are thus pared down during SL analysis, but not sufficiently for the final selection of a TL verb. For example, Levin+ class 9.4 has three possible WordNet senses for *drop*. However, the WordNet sense 8 is not associated with any of the other classes; thus, it is considered to have a higher “information content” than the others. The upshot is that the lexical-selection routine prefers *dejar caer* over other translations such as *derribar* and *bajar*.³ The other classes are similarly associated with ap-

²For further information see the WordNet manuals, section 7, SEMCOR at <http://www.cogsci.princeton.edu>.

³This lexical-selection approach is an adaptation of the notion of *reduction in entropy*, measured by information gain (Mitchell, 1997). Using information content to quantify the “value” of a node in the WordNet hierarchy has also been used for measuring semantic similarity in a taxonomy (Resnik, 1999b). More recently, context-based models of disambiguation have been shown to represent significant improvements over the baseline (Bangalore and Rambow, 2000), (Ratnaparkhi, 2000).

Levin+	Grid/Example	WN Sense	Spanish Verb(s)
9.4 Directional Put	[ag th mod-loc src goal] <i>I dropped the stone</i>	1. move, displace 2. descend, fall, go down 8. drop set down, put down	1. derribar, echar 2. bajar, caerse 8. dejar caer, echar, soltar
45.6 Calibratable Change of State	[th] <i>Prices dropped</i>	1. move, displace 3. decline, go down, wane	1. derribar, echar 3. disminuir
47.7 Meander	[th src goal] <i>The river dropped from the lake to the sea</i>	2. descend, fall, go down 4. sink, drop, drop down	2. bajar, caerse 4. hundir, caer
51.3.1 Roll I	[th goal] <i>The ball dropped to the ground</i>	2. descend, fall, go down	2. bajar, caerse
51.3.1 Roll II	[th particle(down)] <i>The ball dropped down</i>	2. descend, fall, go down	2. bajar, caerse

Table 1: Relation Between Levin+ and WN Senses for ‘drop’

appropriate TL verbs during lexical selection: *disminuir* (class 45.6), *hundir* (class 47.7), and *bajar* (class 51.3.1).⁴

3 Training Data

We began with the lexical database of (Dorr and Jones, 1996), which contains a significant number of WordNet-tagged verb entries. Some of the assignments were in doubt, since class splitting had occurred subsequent to those assignments, with all old WordNet senses carried over to new subclasses. New classes had also been added since the manual tagging. It was determined that the tagging for only 1791 entries—including 1442 verbs in 167 classes—could be considered stable; for these entries, 2756 assignments of WordNet senses had been made. Data for these entries, taken from both WordNet and the verb lexicon, constitute the training data for this study.

The following probabilities were generated from the training data:

- **Grid probability** $r_x = \frac{|{\{r_x \& G_1=G_2\}}|}{|{\{r_x\}}|}$, where r_x is a relation (of relationship type x , e.g., synonymy) between two synsets, s_1 and s_2 , where s_1 is mapped to by a verb in Grid class G_1 and s_2 is mapped to by a verb in Grid class G_2 .

⁴The full set of Spanish translations is selected from WordNet associations developed in the EuroWordNet effort (Dorr et al., 1997).

This is the probability that if one synset is related to another through a particular relationship type, then a verb mapped to the first synset will belong to the same Grid class as a verb mapped to the second synset. Computed values generally range between .3 and .35.

- **Levin + probability** $r_x = \frac{|{\{r_x \& L_{+1}=L_{+2}\}}|}{|{\{r_x\}}|}$, where r_x is as above, except that s_1 is mapped to by a verb in Levin+ class L_{+1} and s_2 is mapped to by a verb in Levin+ class L_{+2} . This is the probability that if one synset is related to another through a particular relationship type, then a verb mapped to the first synset will belong to the same Levin+ class as a verb mapped to the second synset. Computed values generally range between .25 and .3.

- **Tot frame probability** $i,j = \frac{|{\{\theta_{i,v} \& cf_{j,v}\}}|}{|{\{\theta_{i,v}\}}|}$, where $\theta_{i,v}$ is the occurrence of the entire θ -grid i for verb entry v and $cf_{j,v}$ is the occurrence of the entire frame sequence j for a WordNet sense to which verb entry v is mapped. This is the probability that a verb in a Levin+ class is mapped to a WordNet verb sense with some specific combination of frames. Values average only .11, but in some cases the probability is 1.0.

- **Indv frame probability** $i,j = \frac{|{\{\theta_{i,v} \& cf_{j,v}\}}|}{|{\{\theta_{i,v}\}}|}$, where $\theta_{i,v}$ is the occurrence of the single θ -grid component i for verb entry v and $cf_{j,v}$ is the occur-

rence of the single frame j for a WordNet sense to which verb entry v is mapped. This is the probability that a verb in a Levin+ class with a particular θ -grid component (possibly among others) is mapped to a WordNet verb sense assigned a specific frame (possibly among others). Values average .20, but in some cases the probability is 1.0.

• **Prior WN probability** $_s = \frac{|\{t_s\}|}{|\{t_v\}|}$, where t_s is an occurrence of tag s (for a particular synset) in SEMCOR and t_v is an occurrence of any of a set of tags for verb v in SEMCOR, with s being one of the senses possible for verb v . This probability is the prior probability of specific WordNet verb senses. Values average .11, but in some cases the probability is 1.0.

In addition to the foregoing data elements, based on the training set, we also made use of a semantic similarity measure, which reflects the confidence with which a verb, given the total set of verbs assigned to its Levin+ class, is mapped to a specific WordNet sense. This represents an implementation of a class disambiguation algorithm (Resnik, 1999a), modified to run against the WordNet verb hierarchy.⁵

We also made a powerful “same-synset assumption”: If (1) two verbs are assigned to the same Levin+ class, (2) one of the verbs v_1 has been mapped to a specific WordNet sense s_1 , and (3) the other verb v_2 has a WordNet sense s_2 synonymous with s_1 , then v_2 should be mapped to s_2 . Since WordNet groups synonymous word senses into “synsets,” s_1 and s_2 would correspond to the same synset. Since Levin+ verbs are mapped to WordNet senses via their corresponding synset identifiers, when the set of conditions enumerated above are met, the two verb entries would be mapped to the same WordNet synset.

As an example, the two verbs *tag* and *mark* have been assigned to the same Levin+ class. In WordNet, each occurs in five synsets, only one in which they both occur. If *tag* has a WordNet synset assigned to it for the Levin+ class it shares with *mark*, and it is the synset that covers senses

⁵The assumption underlying this measure is that the appropriate word senses for a group of semantically related words should themselves be semantically related. Given WordNet’s hierarchical structure, the semantic similarity between two WordNet senses corresponds to the degree of informativeness of the most specific concept that subsumes them both.

of both *tag* and *mark*, we can safely assume that that synset is also appropriate for *mark*, since in that context, the two verb senses are synonymous.

4 Evaluation

Subsequent to the culling of the training set, several processes were undertaken that resulted in full mapping of entries in the lexical database to WordNet senses. Much, but not all, of this mapping was accomplished manually.

Each entry whose WordNet senses were assigned manually was considered by at least two coders, one coder who was involved in the entire manual assignment process and the other drawn from a handful of coders working independently on different subsets of the verb lexicon. In the manual tagging, if a WordNet sense was considered appropriate for a lexical entry by any one of the coders, it was assigned. Overall, 13452 WordNet sense assignments were made. Of these, 51% were agreed upon by multiple coders. The kappa coefficient (K) of intercoder agreement was .47 for a first round of manual tagging and (only) .24 for a second round of more problematic cases.⁶

While the full tagging of the lexical database may make the automatic tagging task appear superfluous, the low rate of agreement between coders and the automatic nature of some of the tagging suggest there is still room for adjustment of WordNet sense assignments in the verb database. On the one hand, even the higher of the kappa coefficients mentioned above is significantly lower than the standard suggested for good reliability ($K > .8$) or even the level where tentative conclusions may be drawn ($.67 < K < .8$) (Carletta, 1996), (Krippendorff, 1980). On the other hand, if the automatic assignments agree with human coding at levels comparable to the degree of agreement among humans, it may be used to identify current assignments that need review

⁶The kappa statistic measures the degree to which pairwise agreement of coders on a classification task surpasses what would be expected by chance; the standard definition of this coefficient is: $K = (P(A) - P(E)) / (1 - P(E))$, where $P(A)$ is the actual percentage of agreement and $P(E)$ is the expected percentage of agreement, averaged over all pairs of assignments. Several adjustments in the computation of the kappa coefficient were made necessary by the possible assignment of multiple senses for each verb in a Levin+ class, since without prior knowledge of how many senses are to be assigned, there is no basis on which to compute $P(E)$.

and to suggest new assignments for consideration.

In addition, consistency checking is done more easily by machine than by hand. For example, the same-synset assumption is more easily enforced automatically than manually. When this assumption is implemented for the 2756 senses in the training set, another 967 sense assignments are generated, only 131 of which were actually assigned manually. Similarly, when this premise is enforced on the entirety of the lexical database of 13452 assignments, another 5059 sense assignments are generated. If the same-synset assumption is valid and if the senses assigned in the database are accurate, then the human tagging has a recall of no more than 73%.

Because a word sense was assigned even if only one coder judged it to apply, human coding has been treated as having a precision of 100%. However, some of the solo judgments are likely to have been in error. To determine what proportion of such judgments were in reality precision failures, a random sample of 50 WordNet senses selected by only one of the two original coders was investigated further by a team of three judges. In this round, judges rated WordNet senses assigned to verb entries as falling into one of three categories: definitely correct, definitely incorrect, and arguable whether correct. As it turned out, if any one of the judges rated a sense definitely correct, another judge independently judged it definitely correct; this accounts for 31 instances. In 13 instances the assignments were judged definitely incorrect by at least two of the judges. No consensus was reached on the remaining 6 instances. Extrapolating from this sample to the full set of solo judgments in the database leads to an estimate that approximately 1725 (26% of 6636 solo judgments) of those senses are incorrect. This suggests that the precision of the human coding is approximately 87%.

The upper bound for this task, as set by human performance, is thus 73% recall and 87% precision. The lower bound, based on assigning the WordNet sense with the greatest prior probability, is 38% recall and 62% precision.

5 Mapping Strategies

Recent work (Van Halteren et al., 1998) has demonstrated improvement in part-of-speech tag-

ging when the outputs of multiple taggers are combined. When the errors of multiple classifiers are not significantly correlated, the result of combining votes from a set of individual classifiers often outperforms the best result from any single classifier. Using a voting strategy seems especially appropriate here: The measures outlined in Section 3 average only 41% recall on the training set, but the senses picked out by their highest values vary significantly.

The investigations undertaken used both *simple* and *aggregate* voters, combined using various voting strategies. The simple voters were the 7 measures previously introduced.⁷ In addition, three aggregate voters were generated: (1) the product of the simple measures (smoothed so that zero values wouldn't offset all other measures); (2) the weighted sum of the simple measures, with weights representing the percentage of the training set assignments correctly identified by the highest score of the simple probabilities; and (3) the maximum score of the simple measures.

Using these data, two different types of voting schemes were investigated. The schemes differ most significantly on the circumstances under which a voter casts its vote for a WordNet sense, the size of the vote cast by each voter, and the circumstances under which a WordNet sense was selected. We will refer to these two schemes as *Majority Voting Scheme* and *Threshold Voting Scheme*.

5.1 Majority Voting Scheme

Although we do not know in advance how many WordNet senses should be assigned to an entry in the lexical database, we assume that, in general, there is at least one. In line with this intuition, one strategy we investigated was to have both simple and aggregate measures cast a vote for whichever sense(s) of a verb in a Levin+ class received the highest (non-zero) value for that measure. Ten variations are given here:

- **PriorProb:** Prior Probability of WordNet senses
- **SemSim:** Semantic Similarity

⁷Only 6 measures (including the semantic similarity measure) were set out in the earlier section; the measures total 7 because Indv frame probability is used in two different ways.

- **SimpleProd**: Product of all simple measures
- **SimpleWtdSum**: Weighted sum of all simple measures
- **MajSimpleSgl**: Majority vote of all (7) simple voters
- **MajSimplePair**: Majority vote of all (21) pairs of simple voters⁸
- **MajAggr**: Majority vote of SimpleProd and SimpleWtdSum
- **Maj3Best**: Majority vote of SemSim, SimpleProd, and SimpleWtdSum
- **MajSgl+Aggr**: Majority vote of MajSimpleSgl and MajAggr
- **MajPair+Aggr**: Majority vote of MajSimplePair and MajAggr

Table 2 gives recall and precision measures for all variations of this voting scheme, both with and without enforcement of the same-synset assumption. If we use the harmonic mean of recall and precision as a criterion for comparing results, the best voting scheme is MajAggr, with 58% recall and 72% precision without enforcement of the same-synset assumption. Note that if the same-synset assumption is correct, the drop in precision that accompanies its enforcement mostly reflects inconsistencies in human judgments in the training set; the true precision value for MajAggr after enforcing the same-synset assumption is probably close to 67%.

Of the simple voters, only PriorProb and SemSim are individually strong enough to warrant discussion. Although PriorProb was used to establish our lower bound, SemSim proves to be the stronger voter, bested only by MajAggr (the majority vote of SimpleProd and SimpleWtdSum) in voting that enforces the same-synset assumption. Both PriorProb and SemSim provide better results than the majority vote of all 7 simple voters (MajSimpleSgl) and the majority vote of all 21 pairs of simple voters (MajSimplePair). Moreover, the inclusion of MajSimpleSgl and MajSimplePair in a majority vote with MajAggr (in MajSgl+Aggr

⁸A pair cast a vote for a sense if, among all the senses of a verb, a specific sense had the highest value for both measures.

Variation	W/O SS		W/ SS	
	R	P	R	P
PriorProb	38%	62%	45%	46%
SemSim	56%	71%	60%	55%
SimpleProd	51%	74%	57%	55%
SimpleWtdSum	53%	77%	58%	56%
MajSimpleSgl	23%	71%	30%	48%
MajSimplePair	38%	60%	45%	43%
MajAggr	58%	72%	63%	53%
Maj3Best	52%	78%	57%	57%
MajSgl+Aggr	44%	74%	50%	54%
MajPair+Aggr	49%	77%	55%	57%

Table 2: Recall (R) and Precision (P) for Majority Voting Scheme, Before (W/O) and After (W/) Enforcement of the Same-Synset (SS) Assumption

Variation	R	P
AutoMap+	61%	54%
AutoMap-	61%	54%
Triples	63%	52%
Combo	53%	44%
Combo&Auto	59%	45%

Table 3: Recall (R) and Precision (P) for Threshold Voting Scheme

and MapPair+Aggr, respectively) turn in poorer results than MajAggr alone.

The poor performance of MajSimpleSgl and MajSimplePair do not point, however, to a general failure of the principle that multiple voters are better than individual voters. SimpleProd, the product of all simple measures, and SimpleWtdSum, the weighted sum of all simple measures, provide reasonably strong results, and a majority vote of the both of them (MajAggr) gives the best results of all. When they are joined by SemSim in Maj3Best, they continue to provide good results.

The bottom line is that SemSim makes the most significant contribution of any single simple voter, while the product and weighted sums of all simple voters, in concert with each other, provide the best results of all with this voting scheme.

5.2 Threshold Voting Scheme

The second voting strategy first identified, for each simple and aggregate measure, the threshold

value at which the product of recall and precision scores in the training set has the highest value if that threshold is used to select WordNet senses. During the voting, if a WordNet sense has a higher score for a measure than its threshold, the measure votes for the sense; otherwise, it votes against it. The weight of the measure's vote is the precision-recall product at the threshold. This voting strategy has the advantage of taking into account each individual attribute's strength of prediction.

Five variations on this basic voting scheme were investigated. In each, senses were selected if their vote total exceeded a variation-specific threshold. Table 3 summarizes recall and precision for these variations at their optimal vote thresholds.

In the **AutoMap+** variation, Grid and Levin+ probabilities abstain from voting when their values are zero (a common occurrence, because of data sparsity in the training set); the same-synset assumption is automatically implemented. **AutoMap-** differs in that it disregards the Grid and Levin+ probabilities completely. The **Triples** variation places the simple and composite measures into three groups, the three with the highest weights, the three with the lowest weights, and the middle or remaining three. Voting first occurs within the group, and the group's vote is brought forward with a weight equaling the sum of the group members' weights. This variation also adds to the vote total if the sense was assigned in the training data. The **Combo** variation is like Triples, but rather than using the weights and thresholds calculated for the single measures from the training data, this variation calculates weights and thresholds for combinations of two, three, four, five, six, and, seven measures. Finally, the **Combo&Auto** variation adds the same-synset assumption to the previous variation.

Although not evident in Table 3 because of rounding, AutoMap- has slightly higher values for both recall and precision than does AutoMap+, giving it the highest recall-precision product of the threshold voting schemes. This suggests that the Grid and Levin+ probabilities could profitably be dropped from further use.

Of the more exotic voting variations, Triples voting achieved results nearly as good as the AutoMap voting schemes, but the Combo schemes

fell short, indicating that weights and thresholds are better based on single measures than combinations of measures.

6 Conclusions and Future Work

The voting schemes still leave room for improvement, as the best results (58% recall and 72% precision, or, optimistically, 63% recall and 67% precision) fall shy of the upper bound of 73% recall and 87% precision for human coding.⁹ At the same time, these results are far better than the lower bound of 38% recall and 62% precision for the most frequent WordNet sense.

As has been true in many other evaluation studies, the best results come from combining classifiers (MajAggr): not only does this variation use a majority voting scheme, but more importantly, the two voters take into account all of the simple voters, in different ways. The next-best results come from Maj3Best, in which the three best single measures vote. We should note, however, that the single best measure, the semantic similarity measure from SemSim, lags only slightly behind the two best voting schemes.

This research demonstrates that credible word sense disambiguation results can be achieved without recourse to contextual data. Lexical resources enriched with, for example, syntactic information, in which some portion of the resource is hand-mapped to another lexical resource may be rich enough to support such a task. The degree of success achieved here also owes much to the confluence of WordNet's hierarchical structure and SEMCOR tagging, as used in the computation of the semantic similarity measure, on the one hand, and the classified structure of the verb lexicon, which provided the underlying groupings used in that measure, on the other hand. Even where one measure yields good results, several data sources needed to be combined to enable its success.

Acknowledgments

The authors are supported, in part, by PFF/PECASE Award IRI-9629108, DOD

⁹The criteria for the majority voting schemes preclude their assigning more than 2 senses to any single database entry. Controlled relaxation of these criteria may achieve somewhat better results.

Contract MDA904-96-C-1250, DARPA/ITO Contracts N66001-97-C-8540 and N66001-00-28910, and a National Science Foundation Graduate Research Fellowship.

References

- Srinivas Bangalore and Owen Rambow. 2000. Corpus-Based Lexical Choice in Natural Language Generation. In *Proceedings of the ACL*, Hong Kong.
- Olivier Bodenreider and Carol A. Bean. 2001. Relationships among Knowledge Structures: Vocabulary Integration within a Subject Domain. In C.A. Bean and R. Green, editors, *Relationships in the Organization of Knowledge*, pages 81–98. Kluwer, Dordrecht.
- Jean Carletta. 1996. Assessing Agreement on Classification Tasks: The Kappa Statistic. *Computational Linguistics*, 22(2):249–254, June.
- Bonnie J. Dorr and Douglas Jones. 1996. Robust Lexical Acquisition: Word Sense Disambiguation to Increase Recall and Precision. Technical report, University of Maryland, College Park, MD.
- Bonnie J. Dorr and Mari Broman Olsen. 1997. Deriving Verbal and Compositional Lexical Aspect for NLP Applications. In *Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics (ACL-97)*, pages 151–158, Madrid, Spain, July 7–12.
- Bonnie J. Dorr, M. Antonia Martí, and Irene Castellón. 1997. Spanish EuroWordNet and LCS-Based Interlingual MT. In *Proceedings of the Workshop on Interlinguas in MT, MT Summit, New Mexico State University Technical Report MCCS-97-314*, pages 19–32, San Diego, CA, October.
- Christiane Fellbaum. 1998. *WordNet: An Electronic Lexical Database*. MIT Press, Cambridge, MA.
- Eduard Hovy. In press. Comparing Sets of Semantic Relations in Ontologies. In R. Green, C.A. Bean, and S. Myaeng, editors, *The Semantics of Relationships: An Interdisciplinary Perspective*. Book manuscript submitted for review.
- A. Kilgarriff and J. Rosenzweig. 2000. Framework and Results for English SENSEVAL. *Computers and the Humanities*, 34:15–48.
- Klaus Krippendorff. 1980. *Content Analysis: An Introduction to Its Methodology*. Sage, Beverly Hills.
- Beth Levin. 1993. *English Verb Classes and Alternations: A Preliminary Investigation*. University of Chicago Press, Chicago, IL.
- George A. Miller and Christiane Fellbaum. 1991. Semantic Networks of English. In Beth Levin and Steven Pinker, editors, *Lexical and Conceptual Semantics*, pages 197–229. Elsevier Science Publishers, B.V., Amsterdam, The Netherlands.
- Tom Mitchell. 1997. *Machine Learning*. McGraw Hill.
- Mari Broman Olsen, Bonnie J. Dorr, and David J. Clark. 1997. Using WordNet to Posit Hierarchical Structure in Levin’s Verb Classes. In *Proceedings of the Workshop on Interlinguas in MT, MT Summit, New Mexico State University Technical Report MCCS-97-314*, pages 99–110, San Diego, CA, October.
- Martha Palmer. 2000. Consistent Criteria for Sense Distinctions. *Computers and the Humanities*, 34:217–222.
- Adwait Ratnaparkhi. 2000. Trainable methods for surface natural language generation. In *Proceedings of the ANLP-NAACL*, Seattle, WA.
- Philip Resnik. 1999a. Disambiguating noun groupings with respect to wordnet senses. In S. Armstrong, K. Church, P. Isabelle, E. Tzoukermann S. Manzi, and D. Yarowsky, editors, *Natural Language Processing Using Very Large Corpora*, pages 77–98. Kluwer Academic, Dordrecht.
- Philip Resnik. 1999b. Semantic similarity in a taxonomy: An information-based measure and its application to problems of ambiguity in natural language. In *Journal of Artificial Intelligence Research*, number 11, pages 95–130.
- Hans Van Halteren, Jakub Zavrel, and Walter Daelemans. 1998. Improving data-driven wordclass tagging by system combination. In *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and the 17th International Conference on Computational Linguistics*, pages 491–497.