# Automatic Sense Tagging Using Parallel Corpora

Nancy Ide†, Tomaž Erjavec‡ and Dan Tufiş*

†Department of Computer Science
Vassar College
Poughkeepsie, New York 12604-0520, USA
ide@cs.vassar.edu

‡ Department of Intelligent Systems
Institute "Jožef Stefan"
Jamova 39,
SI-1000 Ljubljana, SLOVENIA
tomaz.erjavec@ijs.si

* RACAI
Romanian Academy
CASA ACADEMIEI, Calea 13 Septembrie 13,
Bucharest 74311, ROMANIA
tufis@racai.ro

## Abstract

This article reports the results of an analysis of translation equivalents in six languages from different language families, automatically extracted from an on-line 7-way parallel corpus of George Orwell's *Nineteen Eighty-Four*. The goal is to determine sense distinctions that can be used to automatically sense-tag the data. Our results show that sense distinctions derived from cross-lingual information correspond to those made by human annotators, especially at the coarse-grained level. We also show that the reliability of sense assignments at finer-grained levels is comparable for human annotators and those produced automatically with cross-lingual data.

## 1 Introduction

It is well known that the most nagging issue for word sense disambiguation (WSD) is the definition of just what a word sense is. At its base, the problem is a philosophical and linguistic one that is far from being resolved. However, work in automated language processing has led to efforts to find practical means to distinguish word senses, at least to the degree that they are useful for natural language processing tasks such as summarization, document retrieval, and machine translation. Several criteria have been suggested and exploited to automatically determine the sense of a word in context (see Ide and Véronis, 1998), including syntactic behavior, semantic and pragmatic knowledge, and especially in more recent empirical studies, word co-occurrence within syntactic relations (e.g., Hearst, 1991; Yarowsky, 1993), words co-occurring in global context (e.g., Gale *et al.*, 1993; Yarowsky, 1992; Schütze, 1992, 1993), etc. No clear criteria have emerged, however, and the problem continues to loom large for WSD work.

Resnik and Yarowsky (1997) have recently suggested that for the purposes of WSD, the different senses of a word could be determined by considering only sense distinctions that are lexicalized cross-linguistically. In particular, they propose that some set of target languages be identified, and that the sense distinctions to be considered for language processing applications and evaluation be restricted to those

that are realized lexically in some minimum subset of those languages. This idea would seem to provide an answer, at least in part, to the problem of determining different senses of a word: intuitively, one assumes that if another language lexicalizes a word in two or more ways, there must be a conceptual motivation. If we look at enough languages, we would be likely to find the significant lexical differences that delimit different senses of a word.

Several studies have attempted to use information from parallel texts for WSD (e.g., Gale *et al.,* 1992a, 1993; Dagan *et al.,* 1991; Dagan and Itai, 1994) as well as to define semantic properties of and relations among lexemes (Dyvik, 1998). More recently, two studies have examined the use of cross-lingual lexicalization as a criterion for validating sense distinctions: Ide (1999) used translation equivalents derived from aligned versions of Orwell's *Nineteen Eighty-Four* among five languages from four different languages families, while Resnik and Yarowsky (2000) used translations generated by native speakers presented with isolated sentences in English. In both of these studies, translation information was used to validate sense distinctions provided in lexicons such as WordNet. Although the results are promising, especially for coarse-grained sense distinctions, they rest on the acceptance of a previously established set of senses. Given the substantial divergences among sense distinctions in dictionaries and lexicons, together with the ongoing debate within the WSD community concerning which sense distinctions, if any, are appropriate for language processing applications, fitting cross-linguistic information to pre-established sense inventories may not be the optimal approach.

This paper builds on previously reported work by using translation equivalents derived from parallel corpora to discover sense distinctions that can be used to automatically sense-tag the data. In the next section we explain the methodology used, i.e. the corpus, the automatically derived multilingual lexicon and the subset we used in our experiment, the clustering algorithm and manual assignment of WordNet senses to the testset sentences. Section 3 gives an example of obtained clusters and gives empirical results of comparing the assignment of WordNet senses to those of our program. Section 4 discusses the experiment

and further work, and Section 5 summarises the paper.

## 2   Methodology

We conducted a study using parallel, aligned versions of George Orwell's *Nineteen Eighty-Four[1]* (Erjavec and Ide, 1998, Erjavec 2001) in seven languages: English, Romanian, Slovene, Czech, Bulgarian, Estonian and Hungarian. The study involves languages from four language families (Germanic, Romance, Slavic and Finno-Ugric), three languages from the same family (Czech, Slovene and Bulgarian), as well as two non-Indo-European languages (Estonian and Hungarian).

*Nineteen Eighty-Four* is a text of about 100,000 words, translated directly from the original English to each of the other six languages. The parallel versions of the text are sentence-aligned to the English, where each word is tagged for its context-disambiguated lemma and part of speech, and other word-level morpho-syntactic information, as illustrated in Figure 1.

```
<text id="Oen." lang="en">
<body>
<div type="part" id="Oen.1">
<div type="chapter" id="Oen.1.1">
<p id="Oen.1.1.1">
<s id="Oen.1.1.1.1">
<w lemma="it" ana="Pp3ns">It</w>
<w lemma="be" ana="Vmis3s">was</w>
<w lemma="a" ana="Di">a</w>
<w lemma="bright"
ana="Af">bright</w>
<w lemma="cold" ana="Afp">cold</w>
<w lemma="day" ana="Ncns">day</w>
<w lemma="in" ana="Sp">in</w>
<w lemma="April"
ana="Ncns">April</w>
<c>,</c>
<w lemma="and" ana="Cc-n">and</w>
...
```

**Figure 1. The structure of the corpus**

Although *Nineteen Eighty-Four* is a work of fiction, Orwell's prose is not highly stylized and, as such, it provides a reasonable sample of modern, ordinary language that is not tied to a given topic or sub-domain (which is the case for newspapers, technical reports, etc.).

---

[1] Corpus is available at http://nl.ijs.si/ME/V2/

Furthermore, the translations of the text seem to be relatively faithful to the original: for instance, over 95% of the sentence alignments in the full parallel corpus of seven languages are one-to-one (Priest-Dorman, *et al.,* 1997).

## 2.1 The multilingual lexicon

The first step in the experiment involved automatically constructing a multilingual lexicon, based on the corpus. For this we used the method presented in Tufiş and Barbu (2001), where the translation equivalents of the corpus lemmas are determined on the basis of a 1-1 mapping algorithm that assumes that a given lemma in English is translated by a lemma of the same part of speech appearing sufficiently frequently in the aligned sentences. The desired precision and recall are controlled by various parameters, i.e., log-likelihood scores, string similarity (cognate) scores and occurrence threshold.

Our multilingual lexicon was made in two steps. First, 6 bilingual lexicons with English were automatically extracted from the corpus. Each such lexicon contains about 10,000 entries with an estimated precision and recall of more than 80% each: the log-likelihood threshold was set to 9, minimal cognate score to 0.42 and no threshold for number of occurrences.

From the 6 bilingual lexicons the 7-language multilingual lexicon was then generated, where each English word was associated with all its translations in every language. The resulting size of the multilingual lexicon is 7081 entries, out of which 4042 entries have translation equivalents in all languages.

## 2.2 The lexical dataset

In order to test the validity of our approach we focused the current study on a subset of the multilingual lexicon. We selected only entries where the log likelihood score is at least 18, which have no undetermined translations and at least five ambiguous translations, and have at least ten occurrences in the corpus. From this set we selected the nouns, in order to eliminate variations due to differences of morpho-syntactic category. This yielded a list of 107 English nouns, to each of which the clustering algorithm (described below) was applied. We then further narrowed the list to a sample of 33 words to be sense-tagged by human annotators

for validation and comparison purposes. The final list includes words covering a range of frequencies and degrees of ambiguity (see Appendix).

## 2.3 Clustering

Our method is similar to that of Schütze (1992), who utilized context words to determine sense distinctions. For each noun in the lexical dataset, we extracted all sentences from the English *1984* containing the lemma in question together with the parallel sentences of each of the six translations. The aligned sentences were then automatically scanned to extract translation equivalents. Sentences in which more than one translation equivalent appears were eliminated; this happens in about 5% of the translations.

On the basis of the extracted example sentences, a vector was created for each occurrence, representing all the possible lexical translations into the six parallel versions: if a given word is used to translate that occurrence for a given sentence, the vector contains a 1 in the corresponding position in the vector, and a 0 otherwise. We then fed the occurrence vectors for each ambiguous word to an unweighted pair group algorithm (Sleath and Sokol, 1973),[2] which clusters the vectors by iteratively merging pairs of clusters (initially, each occurrence vector is a cluster) based on the smallest distance between them.

Figure 1 gives an example of the output of the clustering algorithm for "glass". The labels on the leaves of the tree are identifiers for the sentences in the English Orwell for the sentences in which each of the occurrences appears. The output also provides a *minimum distance* value indicating the degree of similarity between each pair of clusters in the graph, where 0 indicates that the clusters are identical. In our data, the highest distance values (typically, but not always, associated with the outermost cluster) range between 2 and 3.

## 2.4 Validation

To validate the clustering results, we compared the groupings of the occurrences obtained by applying the algorithm with sense assignments made by two human annotators on the basis of

---

[2] We used an implementation of the unweighted pair group algorithm developed by Andreas Stolcke.

WordNet.[3] Our original intent was to ask the annotators to group occurrences without reference to an externally defined sense set, in order to produce data as comparable as possible to the clustering algorithm results. However, this task proved to be inordinately difficult and time-consuming for the annotators, and was eventually abandoned.

To normalize the results of the clustering algorithm and the sense assignments made by the human annotators, we represent each annotated sentence for a word with a vector of length $n(n\text{-}1)/2$, where $n$ is the number of occurrences of the word in the corpus. The positions in the vector represents a "yes-no" assignment for each pair of occurrences, indicating whether or not they were judged to belong to the same sense group. Vectors for each word and each annotator were created on the basis of whether or not the paired occurrences had been assigned the same WordNet sense. Representing the clustering algorithm results in this form required some means to "flatten" the cluster hierarchies to conform more closely to the WordNet-based data. To accomplish this, we combined clusters with a minimum distance value at or below 1.7 together, and treated each leaf of the resulting collapsed tree as a different sense. This yielded a set of sense distinctions for each word roughly similar in number to those assigned by the annotators. Note that we used the actual number of senses annotators assigned rather than the number of WordNet senses as a guide to determine the minimum distance cutoff, since it is highly likely that some WordNet senses are not represented in the corpus.

## 3    Results

### 3.1    Clustering

The cluster output for "glass" in Figure 1 is an example of the clustering results we have obtained for the 107 English nouns in this study. The numbers give the distances between the clusters, and the parenthesis contain WordNet senses, manually  assigned for validation; the definitions are given Figure 2. The top group is further divided into two sub-clusters, the lower of which refer to a looking glass and a

magnifying glass, respectively. Interestingly, the clustering in both the top and lower groups reveals additional sub-groupings that are not distinguished in WordNet:  the top sub-group of the top cluster in Figure 1 contains occurrences of "glass"  which deal with some physical aspect of the material ("texture of", "surface of", "rainwatery", "soft", etc.). In the lower cluster, the two main sub-groups distinguish a (drinking) glass as a manipulatable  object (by washing, holding, on a shelf, etc.) from its sense as a vessel (mainly used as the object of "pour into", "fill", "take/pick up", etc. or modified by "empty", "of gin", etc.).

```
1. a brittle transparent solid with
   irregular atomic structure
2. a glass container for holding
   liquids while drinking
3. the quantity a glass will hold
4. a small refracting telescope
5. a mirror; usually a ladies'
   dressing mirror
6. glassware collectively; "She
   collected old glass"
```

Figure 2. WordNet 1.6 senses for *glass*(noun)

### 3.2    Comparison to human annotators

Our results are summarized in Table 1, which gives the percentage of agreement between the cluster algorithm and each annotator, between the two annotators, and for the algorithm and both annotators taken together.  We give here the raw percentages only; common measures of annotator agreement such as the Kappa statistic (Carletta, 1996) proved to be inappropriate for our  two-category  ("yes-no")  classification scheme. The percentages are similar to those reported in earlier work; for example, Ng *et al.* (1999) achieved a raw percentage score of 58% agreement among annotators tagging nouns with WordNet 1.6 senses.

| | |
|---|---|
| Cluster/Annotator 1 | 66.7% |
| Cluster/Annotator 2 | 63.6% |
| Annotator 1/Annotator 2 | 76.3% |
| Cluster/Annotator 1/ Annotator 2 | 53.4% |

Table 1. Levels of agreement

---

[3] Version 1.6,  http://www.cogsci.princeton.edu/~wn

```
          1.452     _____|-> (1) Oen.1.8.100.2
          |-------|         |-> (1) Oen.1.8.101.5
          |       |_____  |-----> (1) Oen.1.8.104.4
          |               | |_____|-> (1) Oen.1.8.65.1
          |               | |      |-> (1) Oen.2.5.14.5
      1.441 |             |-----> (1) Oen.1.8.65.2
      |-------|           |                _|-> (1) Oen.1.8.67.3
      |      |            |              |-| |-> (1) Oen.1.8.70.3
      |      |     |-----||            |-| |-> (1) Oen.2.4.51.8
      |    1.441   |      |          |-| |-> (1) Oen.2.4.66.4
  1.626 |   |-------|     |        |-| |-> (1) Oen.2.4.66.6
  |-------|  |             |-----| |-> (1) Oen.2.7.3.3
  |      |   |             |       |-> (1) Oen.2.9.14.7.5
  |      |   |             |_____|-----> (6) Oen.1.8.63.6
  |      |   |                   |_____|-----> (1) Oen.1.8.9.7
  |      |   |                   |      |-----> (1) Oen.2.10.24.5
  |    1.500  _____|-----> (1) Oen.1.8.72.4
 -| 2.665 |--------|        |-----> (5) Oen.3.3.66.4
  |      |-------> (4) Oen.3.4.7.4
  |    1.558 |-----> (2) Oen.1.8.32.1
  |    |-------|       _|-> (2) Oen.1.8.21.10
  |    |       |     |-| |-> (2) Oen.1.8.25.3
  | 1.599 |    |-----|  |-> (2) Oen.1.8.37.2
  |-------|    |        |-> (2) Oen.1.8.57.5
  |      |     |-------> (2) Oen.2.8.64.1
  | 1.000 |         1.118 _____|-----> (2) Oen.1.7.15.3
  |-------|         |------|     |-----> (2) Oen.1.7.16.6
  |      |          |      |          _|-> (2) Oen.2.8.19.8
  | 1.311 |         |-----|  |-> (2) Oen.2.8.21.1
  |-------|         |        |-> (2) Oen.3.6.4.4
  | 1.414 _____|-> (3) Oen.2.1.19.8
  |------|      |-> (2) Oen.2.8.21.7
  |           _|-> (2) Oen.3.6.2.1
  |------|  |-> (2) Oen.3.6.30.1
  |-> (2) Oen.3.6.39.4
```
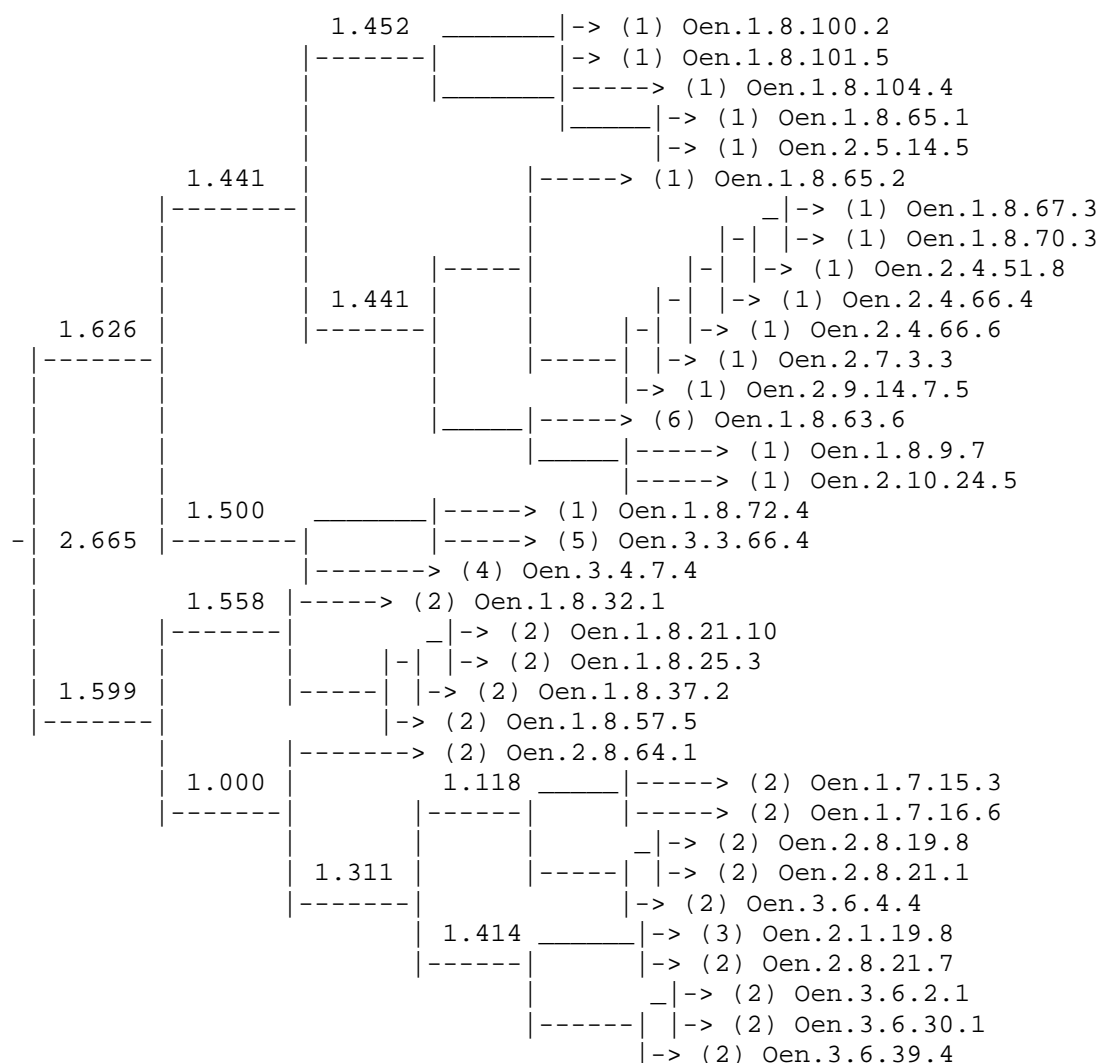
Figure 1. Cluster results for "glass"

## 4    Discussion and further work

Comparison of sense differentiation achieved by considering translation equivalents, as given by the clustering algorithm, with those determined by human annotators suggests that use of translation equivalents for word sense tagging and disambiguation is worth pursuing. Although agreement levels are not astounding, they are comparable to (and in some cases higher than) those obtained in earlier studies tagging with WordNet senses. Furthermore, the difference in agreement between the human annotators and the annotators and the clustering algorithm is only 10-13%, and remains similar to scores obtained in other studies. Given the differences between the nature of the tagging by human annotators, in which pre-defined senses were assigned, and the clustering strategy, the level of agreement is encouraging.

Many studies have pointed out that coarser-grained sense distinctions can be assigned more reliably by human annotators than finer distinctions such as those in WordNet. To address this problem, Ng *et al.* (1999) and Bruce and Wiebe (1998) have recently proposed algorithms that identify coarser-grained distinctions based on tags assigned by

human annotators. In our data, we notice that the clustering algorithm reliably distinguishes coarse-grained senses (e.g., homographs), which is not surprising since translations are far more likely to differ for homographs. Our data suggests that translation equivalents provide a potentially useful means to identify coarse-grained distinctions that are especially relevant for language understanding tasks such as machine translation.

The sense distinctions derived from the clustering algorithm are hierarchical, often identifying four or five levels of refinement, whereas the WordNet sense distinctions are organized as a flat list with no indication of their degree of relatedness. Our attempt to flatten the cluster data in fact loses much information about the relatedness of senses. As a result, annotators (and, in our study, the clustering algorithm) are penalized as much for failing to distinguish senses 2 and 4 of "glass" (given in Figure 2) as for senses 2 and 3 or senses 1 and 6, which are clearly more closely related. We are currently experimenting with utilizing the hierarchy revealed by the cluster data, by devising measures that take relatedness into account. This is in line with the suggestion made by Resnik and Yarowsky (1997) to penalize WSD programs less for failing to distinguish more closely related senses.

## 5    Summary

This study suggests that cross-lingual information can be used for automatic sense tagging that is as reliable as that obtained from human annotators, in particular for relatively coarse-grained distinctions. The clusters derived by our method identify only those occurrences that are more or less closely related—they do not provide a description of the sense such as one would find in a dictionary. While this could be seen as a drawback, it is important to note that WSD studies use *only the knowledge that a set of occurrences of an ambiguous word are used in the same sense*, not the definitions themselves.[4] For example, WSD systems

typically rely on information gathered about the context of occurrences used in the same sense. The "definition" of the word is, in many ways, irrelevant to the exercise.

Our approach is fully automated through all its steps: sentence alignment of the parallel texts, extraction of translation equivalents, identification of translation for each occurrence, and generation of sense clusters that provide the information to tag occurrences deemed to use the same sense of an ambiguous word. The greatest obstacle to its application is the lack of parallel corpora. The freely available parallel corpora for several languages that exist are small (e.g., the Orwell), domain dependent (e.g. the MULTEXT *Journal of the Commission* corpus; Ide and Véronis, 1994) or represent highly stylized language (e.g. the Bible; Resnik *et al.*, 1999). Additional resources will be required to answer larger questions about the use of cross-lingual information for sense tagging, such as the effect of domain and style, and to verify the method using much larger samples.

## References

Bruce, Rebecca and Wiebe, Janyce (1998). Word sense distinguishability and inter-coder agreement. *Proceedings of the Third Conference on Empirical Methods in natural Language Processing.* June, Granada, Spain, 53-60.

Carletta, Jean (1996). Assessing Agreement on Classification Tasks: The Kappa Statistic. *Computational Linguistics*, 22(2), 249-254.

Dagan, Ido and Itai, Alon (1994). Word sense disambiguation using a second language monolingual corpus. *Computational Linguistics*, 20(4), 563-596.

Dagan, Ido; Itai, Alon; and Schwall, Ulrike (1991). Two languages are more informative than one. *Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics,* 18-21 June 1991, Berkeley, California, 130-137.

Dyvik, Helge (1998). Translations as Semantic Mirrors. Proceedings of Workshop W13: Multilinguality in the Lexicon II, The 13th Biennial European Conference on Artificial Intelligence (ECAI 98), Brighton, UK, 24-44.

---

[4] Definitions are sometimes used to provide additional collocates, etc. but are not in themselves required for WSD.

Erjavec, Tomaž and Ide, Nancy (1998). The MULTEXT-EAST Corpus. *Proceedings of the First International Conference on Language Resources and Evaluation*, 27-30 May 1998, Granada, 971-74.

Erjavec, Tomaž (2001). Harmonised Morphosyntactic Tagging for Seven Languages and Orwell's 1984. *Proceedings NLPRS 2001*. Tokyo

Gale, William A., Church, Kenneth W. and Yarowsky, David (1993). A method for disambiguating word senses in a large corpus. *Computers and the Humanities,* 26, 415-439.

Hearst, Marti A. (1991). Noun homograph disambiguation using local context in large corpora. *Proceedings of the 7th Annual Conference of the University of Waterloo Centre for the New OED and Text Research,* Oxford, United Kingdom, 1-19.

Ide, Nancy (1999). Cross-lingual sense determination: Can it work? *Computers and the Humanities,* 34:1-2, 223-34.

Ide, Nancy and Véronis, Jean (1998). Word sense disambiguation: The state of the art. *Computational Linguistics,* 24:1, 1-40.

Ide, N. and Véronis, J. (1994). *Multext (Multilingual Tools and Corpora).* Proceedings of the 14th International Conference on Computational Linguistics, COLING'94, Kyoto, 90-96

Leacock, Claudia; Towell, Geoffrey and Voorhees, Ellen (1993). Corpus-based statistical sense resolution. *Proceedings of the ARPA Human Language Technology Workshop*, San Francisco, Morgan Kaufman.

Melamed, I. Dan. (1997). Measuring Semantic Entropy. *ACL-SIGLEX Workshop Tagging Text with Lexical Semantics: Why, What, and How?* April 4-5, 1997, Washington, D.C., 41-46.

Miller, George A.; Beckwith, Richard T. Fellbaum, Christiane D.; Gross, Derek and Miller, Katherine J. (1990). WordNet: An on-line lexical database. *International Journal of Lexicography,* 3(4), 235-244.

Ng, Hwee Tou, & Lim, Chung Yong, & Foo, Shou King (1999). A Case Study on Inter-Annotator Agreement for Word Sense Disambiguation. *Proceedings of the ACL SIGLEX Workshop on Standardizing Lexical Resources (SIGLEX99),* College Park, Maryland, 9-13.

Priest-Dorman, Greg; Erjavec, Tomaž; Ide, Nancy and Petkevic, Vladimír (1997). Corpus Markup. MULTEXT-East Deliverable D2.3 F.

Resnik, Philip; Broman Olsen, Mari and Diab, Mona (1999). Creating a Parallel Corpus from the Book of 2000 Tongues. *Computers and the Humanities.* Vol. 33: 129-153.

Resnik, Philip and Yarowsky, David. 2000. Distinguishing systems and distinguishing senses: New evaluation methods for word sense disambiguation. *Journal of Natural Language Engineering,* 5(2): 113-133.

Resnik, Philip and Yarowsky, David (1997). A perspective on word sense disambiguation methods and their evaluation. *ACL-SIGLEX Workshop Tagging Text with Lexical Semantics: Why, What, and How?* April 4-5, 1997, Washington, D.C., 79-86.

Schütze, Hinrich (1992). Dimensions of meaning. *Proceedings of Supercomputing'92.* IEEE Computer Society Press, Los Alamitos, California, 787-796.

Schütze, Hinrich (1993). Word space. In Hanson, Stephen J.; Cowan, Jack D.; and Giles, C. Lee (Eds.) *Advances in Neural Information Processing Systems 5*, Morgan Kauffman, San Mateo, California, 5, 895-902.

Sleath, Peter H. A. and Sokal, Robert R. (1973). *Numerical taxonomy; the principles and practice of numerical classification.* W. H. Freeman, San Francisco.

Tufiş, Dan and Barbu, Ana Maria (2001), Automatic Construction of Translation Lexicons. *Proceedings of the WSES and IEEE International Conference on Multimedia, Internet, Video Technologies,* 1-6 September, Malta, 2181-2186.

Yarowsky, David (1992). Word sense disambiguation using statistical models of Roget's categories trained on large corpora. *Proceedings of the 14th International Conference on Computational Linguistics,* COLING'92, 23-28 August, Nantes, France, 454-460.

Yarowsky, David (1993). One sense per collocation. *Proceedings of the ARPA Human Language Technology Workshop,* Princeton, New Jersey, 266-271.

## Appendix A. The Lexical Dataset

Column 1 gives the word, column 2 gives the number of occurrences in the corpus, and the next six columns the number of different words it is translated to according to the automatically derived lexicon, for each of the translation languages. The last three columns give the number of WordNet senses; the average number of multiple translations, i.e. ambiguities; and the average number of WordNet senses actually assigned to occurrences in the corpus.

| Word | No. Occ | RO | SL | CS | BG | ET | HU | WordNet senses | Average number of translations | Average senses assigned |
|---|---|---|---|---|---|---|---|---|---|---|
| act | 35 | 2 | 2 | 2 | 3 | 2 | 2 | 6 | 2.2 | 3 |
| area | 8 | 3 | 3 | 2 | 1 | 2 | 1 | 6 | 2.0 | 2 |
| belief | 14 | 2 | 2 | 2 | 3 | 4 | 3 | 3 | 2.7 | 1.5 |
| bell | 4 | 3 | 2 | 1 | 1 | 3 | 3 | 8 | 2.2 | 3.5 |
| body | 68 | 2 | 1 | 2 | 3 | 2 | 2 | 9 | 2.0 | 2.5 |
| book | 49 | 2 | 2 | 2 | 3 | 2 | 1 | 8 | 2.0 | 1 |
| boot | 15 | 1 | 3 | 3 | 2 | 3 | 3 | 4 | 2.5 | 2 |
| boy | 16 | 3 | 2 | 3 | 2 | 2 | 4 | 4 | 2.7 | 2 |
| breast | 9 | 2 | 1 | 6 | 2 | 1 | 2 | 3 | 2.3 | 2.5 |
| cent | 6 | 3 | 2 | 2 | 2 | 1 | 2 | 2 | 2.0 | 2 |
| condition | 6 | 2 | 2 | 1 | 2 | 2 | 3 | 7 | 2.0 | 4.5 |
| contact | 11 | 3 | 3 | 2 | 3 | 2 | 3 | 8 | 2.7 | 5.5 |
| corner | 24 | 2 | 3 | 3 | 3 | 4 | 3 | 11 | 3.0 | 5 |
| country | 18 | 3 | 2 | 3 | 1 | 3 | 2 | 5 | 2.3 | 4 |
| date | 11 | 3 | 2 | 4 | 2 | 4 | 2 | 8 | 2.8 | 5 |
| day | 80 | 2 | 2 | 2 | 3 | 2 | 4 | 9 | 2.5 | 1.5 |
| department | 12 | 2 | 2 | 1 | 1 | 5 | 1 | 3 | 2.0 | 2 |
| destruction | 9 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 2.0 | 1.5 |
| dial | 9 | 2 | 3 | 2 | 1 | 3 | 1 | 4 | 2.0 | 4.5 |
| distance | 16 | 2 | 2 | 2 | 2 | 2 | 3 | 6 | 2.2 | 2.5 |
| eye | 72 | 3 | 2 | 4 | 2 | 2 | 5 | 5 | 3.0 | 5.5 |
| face | 113 | 8 | 2 | 7 | 4 | 7 | 5 | 14 | 5.5 | 3 |
| fact | 51 | 3 | 2 | 3 | 2 | 2 | 2 | 4 | 2.3 | 4 |
| feeling | 47 | 2 | 2 | 2 | 1 | 1 | 4 | 7 | 2.0 | 5 |
| figure | 25 | 3 | 3 | 2 | 2 | 2 | 4 | 13 | 2.7 | 1 |
| finger | 36 | 2 | 2 | 2 | 2 | 4 | 2 | 2 | 2.3 | 2.5 |
| girl | 51 | 1 | 4 | 4 | 3 | 2 | 2 | 5 | 2.7 | 5 |
| glass | 35 | 2 | 2 | 4 | 3 | 3 | 2 | 6 | 2.7 | 4 |
| good | 10 | 4 | 1 | 2 | 2 | 3 | 1 | 3 | 2.2 | 6 |
| hand | 96 | 5 | 2 | 4 | 4 | 4 | 3 | 14 | 3.7 | 2 |
| hour | 45 | 2 | 1 | 2 | 3 | 4 | 3 | 4 | 2.5 | 9.5 |
| line | 21 | 4 | 5 | 5 | 2 | 3 | 4 | 29 | 3.8 | 4 |
| movement | 40 | 1 | 2 | 2 | 2 | 3 | 4 | 10 | 2.3 | 3 |