

# Bilingual terminology extraction: an approach based on a multilingual thesaurus applicable to comparable corpora

Hervé Déjean, Éric Gaussier

Xerox Research Centre Europe

6, Chemin de Maupertuis, 38240 Meylan, France

firstname.lastname@xrce.xerox.com

Fatia Sadat

Graduate School of Information Science

Nara Institute of Science and Technology, Nara, Japan

[fatia-s@is.aist-nara.ac.jp](mailto:fatia-s@is.aist-nara.ac.jp)

## Abstract

This paper presents several methods for exploiting multiple resources in bilingual lexicon extraction, either from parallel or comparable corpora. First, a special attention is given to the use of multilingual thesauri, and different search strategies based on such thesauri are investigated. Then, a method to optimally combine the different resources for bilingual lexicon extraction is presented. Our results show that the combination of the resources significantly improves results, and that the use of the hierarchical information contained in our thesaurus, UMLS/MeSH, is of primary importance. Lastly, methods for bilingual terminology extraction and thesaurus enrichment are discussed.

## Introduction

The growing availability of comparable corpora, through the Internet or via distribution agencies providing newspapers articles in different languages, has led researchers to develop methods to extract bilingual lexicons from such corpora, in order to enrich existing bilingual dictionaries and thesauri, and help cross the language barrier for cross-language information retrieval. The results obtained thus far on comparable corpora, even though encouraging, are not completely satisfactory yet. We thus investigate in this paper, whether combining different resources can help in the process of extracting bilingual lexicons from corpora, and propose different methods for exploiting multilingual thesauri.

We first present the different methods we use to extract probabilistic translation lexicons: from a comparable corpus, and from a multilingual thesaurus. We then discuss the problem of integrating different resources, namely a general bilingual dictionary and a specialized multilingual thesaurus (the Medical Subject Headings, MeSH, provided through the metathesaurus Unified

Medical Language System, UMLS). Section 6 presents the experiments we conducted to test our different methods, and, finally, section 7 introduces the method we retained for terminology extraction and thesaurus enrichment. All our examples and experiments are based on the (German,English) language pair.

## Context vectors: a basic building block

Bilingual lexicon extraction from non-parallel but comparable corpora has been studied by a number of researchers, (Rapp, 1995; Peters, 1995; Tanaka, 1996; Shahzad 1999; Fung, 2000) among others. Their work relies on the assumption that if two words are mutual translations, then their more frequent collocates (taken here in a very broad sense) are likely to be mutual translations as well. Based on this assumption, a standard approach consists in building context vectors, for each source and target word, which aim at capturing the most significant collocates. The target context vectors are then translated using a general bilingual dictionary, and compared with the source context vectors. This approach is reminiscent of the way similarities between terms are built in information retrieval, through the use of the cosine measure between term vectors extracted from the term-document matrix (Salton *et al.*, 1983).

The use of a general bilingual dictionary to translate context vectors is justified by the fact that if the context vectors are sufficiently large, then some of their elements are likely to belong to the general language and to the bilingual dictionary, and we can thus expect the translated context vector of word  $t$  to be, in average, closer to the context vector of the translation  $s$  of  $t$ . It has to be noted that the above strategy makes sense even when  $t$  is present in the bilingual dictionary, since the corpus may display a particular, technical usage of  $t$ . Our implementation of this strategy relies on the following steps:

1. for each word  $w$ , build a context vector by considering all the words occurring in a window encompassing several sentences that is run through the corpus. Each word  $i$  in the context vector of  $w$  is then weighted with a measure of its association with  $w$ . We chose the log-

- likelihood ratio test, (Dunning, 1993), to measure this association
2. the context vectors of the target words are then translated with our general bilingual dictionary, leaving the weights unchanged (when several translations are proposed by the dictionary, we consider all of them with the same weight)
  3. the similarity of each source word  $s$ , for each target word  $t$ , is computed on the basis of the cosine measure
  4. the similarities are then normalized to yield a probabilistic translation lexicon,  $P(t/s)$ .

To illustrate the above steps, we give below the first 5 words of the context vector of the German word *Leber* (*liver*), together with their associated score, and the translated vector.

German vector	Translated vector
transplantation 138.897	transplant 38.897
resektion 53.501	tumour 48.654
metastase 41.668	secondary 42.552
arterie 38.519	metastasis 41.668
cirrhose 26.302	artery 38.519

One can note that the German term *Resektion* was not found in our bilingual dictionary, and thus not translated. However, the translated context vector contains English terms characteristic of the co-occurrence pattern for *liver*, allowing one to associate the two words *Leber* and *liver*.

### Lexical translation models from a multilingual thesaurus

A multilingual thesaurus bridges several languages through cross-language correspondences between concept classes (a concept class in the thesaurus links alternative names and views of the same concept together, as *cancer of spleen*, *splenic cancer*, *spleen neoplasms* for *splenic neoplasms*). In many cases, it corresponds to a synonym class). The correspondence can be one-to-one, i.e. the same concept classes are used in the different languages, or many-to-many, i.e. different concept classes are used in different languages, and a given concept class in a given language corresponds to zero, one or more concept classes in the other languages. The correspondence between concept classes across languages helps us write the probability  $P(t/s)$  of selecting word  $t$  as a translation of word  $s$  in the following general way, where  $C$  represents a multilingual concept class in MeSH (we omit the derivation, which is mainly technical, and uses the fact that the

correspondence between concept classes in MeSH is one-to-one):

$$P(t/s) = \sum_C P(C/s) P(t/C, s)$$

a formula which can be interpreted as follows: from a source word  $s$ , select a (interlingual) concept class in the thesaurus, according to  $P(C/s)$ , then generate a target word  $t$  from the concept class and the source word, according to  $P(t/C, s)$ . Furthermore, in order not to privilege any of the possible lexicalizations of a given concept (which can be done via the dependence on  $s$  in the last probability distribution), we make the additional simplifying assumption that, given a concept class, the target word  $t$  is independent of the source word  $s$ , which leads to the simplified formula:

$$P(t/s) \approx \sum_C P(C/s) P(t/C)$$

The above equation views the thesaurus as a trellis linking source and target words. As such, given probabilities  $P(C/s)$  and  $P(t/C)$  (see below for the way we estimate these probabilities), there are several ways to compute an association score between source and target words. The most obvious one is to carry the sum over all concept classes, or a large subset of them, as indicated by the formula. We refer to this method as the **complete search**. However, if the relation between a word and a concept class is not significant, the complete search has the disadvantage of bringing noisy data in the estimation of  $P(t/s)$ . An alternate solution is to select just the concept class which maximizes the association between  $s$  and  $t$ . Because of its analogy with the Viterbi algorithm, we refer to this method as the **Viterbi search**.

Nevertheless, neither the complete nor the Viterbi search makes use of the hierarchical information contained in the thesaurus, which is, in the above formulations, mainly viewed as a specialized lexicon. We present below a third search strategy which directly makes use of the structure of the thesaurus. For reasons that will become clear, we call this strategy the **subtree search**.

### The subtree search

Complete search and the Viterbi search represent two extreme ways of making use of the thesaurus since they consider either all or only one of the concept classes it contains. In order to find a way in-between and to focus on a subset of interesting concept classes, we first select for each source word  $s$  the  $n$  best concept classes in the thesaurus, i.e. the first  $n$  concept classes according to the probability distribution  $P(C/s)$ . We then extend this

set of classes by adding new classes using the hierarchy in the thesaurus.

Intuitively, if two or more classes in the selected subset have the same parent class, then the source word is likely to be related to this parent as well as to the classes themselves, since the parent is the direct node "conceptually" linking the classes. For example, if a source word  $s$  selects the two classes *Hepatitis* and *Cirrhosis*, then  $s$  is likely to be related to *Liver Diseases*, the parent class. We make use of this intuition in the following way: for each pair of classes from the set of the  $n$  best classes associated with source word  $s$ , select the subtree formed by the classes, their common ancestor, and all the nodes that appear between the classes and their ancestor.

This algorithm provides a set of subtrees from the 15 sub-thesauri corresponding to the 15 main categories of the MESH classification (MeSH, rather than being a single thesaurus, contains 15 different sub-thesauri, artificially related through a common root node in UMLS. We do not make use of this distinction in the complete and Viterbi methods, but use it for the subtree search to avoid linking classes via the artificial root concept). One can also note that the above algorithm suggests a way to identify polysemous words, or words used through different points of view, via the different sub-thesauri they select subtrees from. This refinement, which should lead to more fine-grained bilingual lexicons, will be the focus of future research.

The set of classes contained in the subtrees is then used in equation (2) to derive associations between source and target words. Table 1 shows the different behaviors of the complete method, using the 200 concept classes closest to the source word, and the subtree method with  $n = 20$  on the source word *Leber*. Even if some candidates are not actual translations of *Leber*, the subtree search provides *liver* (the correct translation), as the first translation candidate. Note that the subtree search, with  $n = 20$ , yields in average 4 subtrees per source word.

Subtree	Complete
<b>liver</b>	hepatocyte
orthotopic	neoplasm
hepatic	enormous
survival	orthotopic
metastasis	inherit

Table 1: First 5 candidates for *Leber*

## Linking words and concept classes

The estimation of the probability distributions  $P(C/s)$  and  $P(t/C)$  used in equation (2) can be easily carried out by resorting once again to context vectors. Indeed, if a word of the corpus is similar to a term present in a concept class, then they are likely to share similar contexts and have similar context vectors. We thus extend the notion of context vectors to concept classes, and rely again on the cosine measure to compute similarities between words and concept classes. The probability distributions  $P(C/s)$  and  $P(t/C)$  are finally derived through normalization.

To build a context vector for a concept class, we first build the context vector of each term the class contains. For single-word units, we directly rely on the context vectors extracted in section 2. If the term is a multi-word unit, as *liver disease*, we consider the conjunction of the context vectors of each word in the unit, normalizing the weights by the number of words in the unit. For example, the context vector for *liver disease* will contain only those words that appear in the context of both *liver* and *disease*, since the whole unit is a narrower concept than its constituents. We then take the disjunction of all context vectors of each entry term in the class, normalizing the weights by the number of terms in the class, to build the context vector of each concept class.

The following example illustrates the complete process: the German graphical variant *Actinomykose* is used in our corpus in addition to *Aktinomykose*, which is the only form listed in the UMLS class C0001261; nevertheless, our process associates C0001261 as the closest class to *Actinomykose* and *actinomycosis* (English), and retain them as translation candidates.

## Combining different models

The previous sections provide us with two different probabilistic lexical translation models: one derived from the corpus, and one from the multilingual thesaurus. A third lexical translation model can be obtained from the bilingual dictionary by considering the different translations of a given entry as equiprobable. For example, our dictionary associates *abbilden* with the two words *depict* and *portray*, thus  $P(\text{depict}|\text{abbilden}) = P(\text{portray}|\text{abbilden}) = 0.5$ . Note that these three models are not independent of each other, since the corpus is used, through the estimation of  $P(C/s)$  and  $P(t/C)$ , in the thesaurus-based model, and the bilingual dictionary is used for translating context vectors in the corpus-based model.

The combination of the different models is then realized in the following mixture of models:

$$P(t/s) = \sum_i P(i) P_i(t/s)$$

where  $i$  is an integer used to index the different models (here  $1 \leq i \leq 3$ ), and  $P(i)$  denotes the probability of selecting model  $i$ . We will not enter into the technical presentation of how we estimate mixture weights here, but simply mention that we have to rely on a bilingual lexicon manually extracted from our corpus. Table 2 presents the mixture weights we obtained for 3 different ways of exploiting the thesaurus (Viterbi search, complete search with 200 classes, subtree search with 20 classes).

	Viterbi	Complete	Subtree
corpus	0.59	0.45	0.33
thesaurus	0.1	0.24	0.37
dictionary	0.31	0.31	0.29

Table 2: Mixture weights for the 3 models

As one can note, these results suggest that the thesaurus is less reliable than the other resources when used with the Viterbi search, whereas it becomes the most important resource with the subtree search. This is not surprising if one relates these results with the way the thesaurus is exploited in each case: precise but incomplete information is used in Viterbi search, more complete but less precise information in the complete search, and more complete and more precise information in the subtree search, since only accurate concept classes are selected and expanded through the thesaurus hierarchy.

Lastly, in the context of cross-language information retrieval, data for estimating mixture weights can be automatically derived by building pseudo-translation pairs in which the source words are extracted from the set of queries used in the training phase, and their corresponding target words are taken to be those occurring in all the documents judged to be relevant to the query. Our first experiments in this context showed that the combination thus obtained was always better than a simple merge of the resources, but we nevertheless failed to report significant improvements on the retrieval results. However, these experiments were solely based on a corpus and a bilingual dictionary, since no multilingual thesaurus was available for the domain under consideration. As we have seen above, the thesaurus brings valuable information, and we

believe we should benefit from using one in this context.

## Linguistic preprocessing

As a preprocessing step, we tag and lemmatize texts in both languages. This step allows us to focus on content words only (nouns, verbs, adjectives and adverbs), and reduces the noise in our model (content words are the primary focus for thesaurus enrichment and cross-language information retrieval). Nevertheless, since we use the (German, English) language pair for all our experiments, a major problem still resides in the difference in the word definition between the two languages, mainly due to the particular usage of compounding the German language has. Two alternatives are offered: either use a direct phrasal alignment, or decompose the German compounds into smaller units. Inasmuch as the models presented in the preceding sections implicitly assume a one-to-one correspondence between words in the two languages, we rely on the second strategy.

However, an additional complication is introduced by the fact that our corpora belongs to the medical domain, thus leaving our German lemmatizer clueless when it comes to decomposing medical compounds. We thus used two additional heuristics, recursively applied on all German words:

1. some sequences, e.g. -ungs-, -heits-, -keits-, -schafts-, -aets- and -ions-, as well as their plural forms, are considered as boundaries between two words in a compound, and break a word into two parts
2. if a word is composed of the sequence AB, and if A and B are both longer than 3 characters and both occur in the corpus, then the sequence AB is decomposed into A and B.

The above heuristics reduce the number of different lemmas in the German vocabulary by 28% (from 14,700 to 10,500), while not hurting too much the quality of the vocabulary since their precision is estimated to be above 90%. For example, they allow us to accurately decompose the compound *Adhaesionsileusbehandlung* into the three parts *Adhaesion*, *Ileus* and *Behandlung*.

## Experiments and results

To test the above models and their combination, we used roughly 700 abstracts from MEDLINE<sup>1</sup>, in German and English (each portion, German and English, contains approximately 100,000 words). These abstracts are “partial” translations of each other,

<sup>1</sup> <http://www4.ncbi.nlm.nih.gov/PubMed/>

because in some cases the English writer directly summarizes the articles in English, rather than translating the German abstracts. That set of abstracts is used both as our comparable corpus as our parallel corpus (cf the following section). We are fully aware that our comparable corpus is thus “ideal” since it is close to a parallel one. Unfortunately, we do not have a “fully” comparable corpus in the medical domain that we could use in conjunction with MeSH. Note however that the choice of this corpus does not, as far as we can judge, bias our results when evaluating the different thesaurus-based models and the model combination<sup>2</sup>.

As already mentioned, we manually extracted a reference lexicon comprising 1,800 translation pairs from our comparable corpus. From this, we reserve approximately 1,200 pairs for estimating the mixture weights, and 600 for the evaluation proper. All our results are averaged over 10 different such splits. Since the models we rely on yield a ranked set of translation candidates for each source word, and since one cannot expect the right translation to be *the* first candidate, we compute precision and recall of each method in the following way: for each pair  $(s,t)$  in the evaluation lexicon, we consider the first  $p$  candidates provided for  $s$  by the method under evaluation, and judge the set as correct if it contains  $t$ , as incorrect otherwise; precision is then obtained by dividing the number of correct sets by the number of sets proposed by the method for the words in the evaluation lexicon, whereas recall is obtained by dividing the number of correct sets by the number of pairs in the evaluation lexicon. In addition, we evaluate the average rank of the first correct translation in the proposed list of translations, for each method.

Model	F1-score
Dictionary	56.16
Corpus	62.04
Thesaurus (ST50)	51.34

Table 3: Results for separate models

Table 3 shows the results we obtained on our comparable corpora, for  $p=10$ , without combining

<sup>2</sup> Such a bias is reported in (Masuichi *et al.*, 2000), their method extracting a parallel corpus from a comparable one.

the different models. ST50 refers to the subtree search strategy within the thesaurus, with  $n=50$ . The precision of the dictionary-based model is around 78%, which is not that bad considering the domain we focused on, but, as one can expect, its recall reaches only 48%. The F1-score, which combines precision and recall, obtained for the corpus-based model is similar to the ones obtained in previous works.

	p=5	p=10
Viter	71.3/14.7	79.7/14.7
Comp(100)	75.4/14.1	80.3/14.1
Comp(200)	75.4/12.3	83.2/12.3
ST10	75.8/11	82.4/11
ST20	76.4/11.7	84.1/11.7
ST50	77.3/11.2	83.6/11.2
ST100	76.9/11.8	83/11.8

Table 4: Evaluation of search strategies

Table 4 presents the results (F1-score) we obtained with the different search strategies for the thesaurus-based model (the Viterbi search, the complete one (considering the first 100 and first 200 concept classes for each source word), and the subtree search with different values of  $n$ ), and two different values for  $p$ , 5 and 10. The average rank is given next to each F1-score.

As one can see, the combination significantly improves the results over the models alone, since the F1-score goes from 62% to 84%, a score that may be good enough to consider manual revisions. Furthermore, the best results are obtained with the subtree search, with  $n=20$ , thus validating our hypothesis that using the structure of the thesaurus is beneficial. One can note however that the results obtained with the complete search using 200 classes are close to the best results. Nevertheless, the optimal subtree search (ST20) uses 7.5 times less classes than the complete search, and is also two times faster. This proves that the subtree search is able to focus on accurate concept classes in the thesaurus, whereas the complete search needs considering more classes to reach a comparable level of performance. Interestingly, it also seems that the candidates provided by the subtree search closely correspond to a semantic field, whereas the ones given by the complete search are more varied. Where this to be the case, the subtree search would also certainly outperform the other methods when used for cross-language information retrieval. We will try to validate this hypothesis in future work.

## Bilingual terminology extraction

Bilingual terminology extraction is based on three steps: word alignment, term extraction and term alignment

In this section, we rely on the word to word translation lexicon obtained from the parallel corpus, following the method described in (Gaussier *et al.*, 2000).

### Term extraction

For identifying German and English candidate terms we use the following patterns, similar to those proposed by (Heid, ) and (Blanck, 2000):

1. single words which appear in the thesaurus (for alignment purposes) or which contain English morphemes extracted from *The Specialized Lexicon* found in UMLS and translated in German (ectomy/ektomie)
2. syntactic patterns: [(ADJ)+ NOUN GEN+] and [ADJ+ NOUN (GEN)+] for German, and all non recursive noun phrases for English

### Term alignment

Our algorithm allows alignment of a sequence of candidate terms, and follows the one proposed in (Hull, 1997). We first try to align candidate terms, and then test if a longer unit, composed of several candidate terms, improve the alignment score. A unit is extended if and only if the next contiguous candidate term is a prepositional phrase, the relaxation of this constraint introducing too much noise. The extension stops when the score is lower than the score of the “non-extended term”. For instance, an alignment score is computed for [*problematischen Gebieten der Chirurgie*] and [*problematic fields*]. Then the English term is extended to [[*problematic fields*] of [*surgery*]], which provides a better alignment score, and is then kept. In this particular example, neither the German nor the English units can be further extended, since the German term occurs at the end of a sentence and the English unit is not followed by a prepositional phrase. The German candidate *problematischen Gebieten der Chirurgie* is thus finally aligned with the English candidate *problematic fields of surgery*.

Most German compounds, decomposed for word alignment purposes, are aligned with English terms corresponding to a sequence *adjective+noun* (*Nierenfunktion/renal function*) or *noun+of+noun* (*Lebensqualitaet/quality of live*). Correspondences

between acronyms and translated developed forms can also be found (*Nierenzellcarcinom/RCC*). In practice, no unit composed of three candidate terms is found. The longest units are generated by German candidate term with a genitive structure (*Plattenepithelcarcinom des Oesophagus/squamous cell esophageal cancer*).

	precision	recall
1	56.52	50.98
2	71.01	64.05
5	84.78	76.47
10	89.85	81.04

Table 5: Evaluation of term alignment

We manually extracted 150 candidate terms with their translation for evaluating our procedure.

Table 5 shows precision and recall for our method. Precision is always higher than recall, which can be explained by the fact that the reference terms were extracted manually when the automatic extraction can propose incorrect units due to chunking errors.

### Thesaurus enrichment

We propose in this section some solutions for enriching monolingual as well as bilingual thesauri. Our goal is to propose tools to the terminologist, in order to enrich semi-automatically existing thesauri. Since the English language is predominant in UMLS, enrichment in other languages via term alignment is possible. The German thesaurus we use, DMD, is indeed a partial German translation of MeSH. Our bilingual extracted lexicon thus provides us with a way to extend the German thesaurus.

The first extension is the introduction of new strings (following UMLS terminology) associated with a concept in the thesaurus. If one element of the bilingual extracted lexicon is in the thesaurus, the translated candidate can be directly added in the part of the thesaurus corresponding to its language. These new strings correspond to synonyms as well as spelling or term variants. The German string *Karzinom* is associated with the UMLS concept C0007097. The English string is *carcinoma*. Through the term alignment, we can provide a new German string for this concept: *Carcinom*. Note that this spelling difference is in fact due to two different German spelling used in medical texts. New strings due to morpho-syntactic variations can also be inserted. Thus, our alignment provides a new string for the entry *Lebertransplantation: Transplantation der Leber*, although this entry already contains four different strings.

A second kind of enrichment is the addition of new concepts in one language. In some cases no German string is proposed for a given concept class. For example, the German thesaurus has no associated strings with C0334281 (*malignant insulinoma*) and C0406864 (*flap loss*). The alignment allows us to propose the following candidates: *malignen Insulinom* for C0334281, and *Lappenverlust* for C0406864, that a terminologist can review before entering in the thesaurus.

The most difficult situation is the classification of new concepts in both languages. One solution to this problem is to use the similarity between words and concepts in order to compute the similarity between terms and concepts. A simple way is to sum the similarity of all the words of a term for a given concept. This gives good results when the new term is a narrower term of an existing concept. For instance, no concept exists in UMLS for the German string *chronische Pankreatitis* or its English translation, *chronic pancreatitis*. Computing a similarity between concepts and these terms in the way described above yields C0030305 (*pancreatitis*) as the closest concept class to both candidate terms. A refinement of the above procedure can be obtained by using syntactic filters. Candidates matching the pattern *ADJ NOUN* are more likely to be interesting than other ones. In general, if the term is composed with some words present in the thesaurus, the list of concepts proposed to the terminologist is relevant. Nevertheless, one cannot go to far in this direction, since the choice of the final concept has to be made by the terminologist: the second best concept proposed for *chronic pancreatitis* is for example *chronic disease*, and the choice of the correct concept completely depends on the organization of the thesaurus. In a preceding version of MeSH, *chronic hepatitis B* was linked to *chronic disease* and to *hepatitis B*, a link which is deleted in the current version (2001). Nevertheless, if the candidate is entirely new (no word of it are present in the thesaurus), the above strategy does not make much sense.

A solution in this case is to combine hierarchical information with some particular morpho-syntactic patterns. For example, words with the suffix *-ectomy* occur below the concept C0543467 (*Surgical Procedures*) and words with the suffix *-graphy* occur below *Diagnosis*. Through this

morpho-syntactic information, a part of the thesaurus can be automatically proposed to the terminologist.

## Conclusion

We have shown how to optimally combine different models derived from different resources for bilingual lexicon extraction from comparable corpora. Such a combination significantly (by 30%) improves the results over the models alone. We have also presented different models based on a multilingual thesaurus, and have obtained the best results with the model integrating hierarchical information. Lastly we have proposed different ways to enrich existing thesauri with new terms discovered in corpora.

## References

- Blank, I., 2000. Terminology extraction from parallel technical texts. In J. Véronis (Ed.), *Parallel Text Processing - Alignment and Use of Translation Corpora*.
- Dunning, T., 1993. Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, 19(1):64-74.
- Fung, P., 2000. A statistical view on bilingual lexicon extraction - From parallel corpora to non-parallel corpora. In J. Véronis (Ed.), *Parallel Text Processing - Alignment and Use of Translation Corpora*.
- Gaussier, E., Hull, D., Ait-Mokhtar, S., 2000. Term alignment in use: Machine-aided human translation. In J. Véronis (Ed.), *Parallel Text Processing - Alignment and Use of Translation Corpora*.
- Heid, U., 1999. A linguistic bootstrapping approach to the extraction of term candidates from German text. *Terminology*, 5(2).
- Hull, D., 1997. Automating the construction of bilingual terminology lexicons. *Terminology*, 4(2).
- Masuichi, H., Flournoy, R., Kaufmann, S., Peters, S., 2000. A bootstrapping method for extracting bilingual text pairs. *COLING Proceedings*.
- Peters, C., Picchi, E., 1995. Capturing the comparable: A system for querying comparable text corpora. *JADT Proceedings*.
- Rapp, R., 1995. Identifying word translations in nonparallel texts. *ACL Proceedings*.
- Salton, G., McGill, J., 1983. *Introduction to Modern Information Retrieval*. New York: McGraw-Hill.
- Shahzad, I., Ohtake, K., Masuyama, S., Yamamoto, K., 1999. Identifying translations of compound nouns using non-aligned corpora. *Workshop MAL Proceedings*.
- Tanaka, K., Iwasaki, H., 1996. Extraction of lexical translations from non-aligned corpora. *COLING Proceedings*.
- Vivaldi, J., Rodriguez, H., 2001. Improving term extraction by combining different techniques. *Terminology*, 7(1).