

Mixed Language Query Disambiguation

Pascale FUNG, LIU Xiaohu and CHEUNG Chi Shun

HKUST

Human Language Technology Center

Department of Electrical and Electronic Engineering

University of Science and Technology, HKUST

Clear Water Bay, Hong Kong

{pascale,lxiaohu,eepercyc}@ee.ust.hk

Abstract

We propose a mixed language query disambiguation approach by using co-occurrence information from monolingual data only. A mixed language query consists of words in a *primary language* and a *secondary language*. Our method translates the query into monolingual queries in either language. Two novel features for disambiguation, namely contextual word voting and 1-best contextual word, are introduced and compared to a baseline feature, the nearest neighbor. Average query translation accuracy for the two features are 81.37% and 83.72%, compared to the baseline accuracy of 75.50%.

1 Introduction

Online information retrieval is now prevalent because of the ubiquitous World Wide Web. The Web is also a powerful platform for another application—interactive spoken language query systems. Traditionally, such systems were implemented on stand-alone kiosks. Now we can easily use the Web as a platform. Information such as airline schedules, movie reservation, car trading, etc., can all be included in HTML files, to be accessed by a generic spoken interface to the Web browser (Zue, 1995; DiDio, 1997; Raymond, 1997; Fung et al., 1998a). Our team has built a multilingual spoken language interface to the Web, named SALSA (Fung et al., 1998b; Fung et al., 1998a; Ma and Fung, 1998). Users can use speech to surf the net via various links as well as issue search commands such as “*Show me the latest movie of Jacky Chan*”. The system recognizes commands and queries in English, Mandarin and Cantonese, as well as mixed language sentences.

Until recently, most of the search engines handle keyword based queries where the user types

in a series of strings without syntactic structure. The choice of key words in this case determines the success rate of the search. In many situations, the key words are ambiguous.

To resolve ambiguity, query expansion is usually employed to look for additional keywords. We believe that a more useful search engine should allow the user to input natural language sentences. Sentence-based queries are useful because (1) they are more natural to the user and (2) more importantly, they provide more contextual information which are important for query understanding. To date, the few sentence-based search engines do not seem to take advantage of context information in the query, but merely extracting key words from the query sentence (AskJeeves, 1998; ElectricMonk, 1998).

In addition to the need for better query understanding methods for a large variety of domains, it has also become important to handle queries in different languages. **Cross-language information retrieval** has emerged as an important area as the amount of non-English material is ever increasing (Oard, 1997; Grefenstette, 1998; Ballesteros and Croft, 1998; Picchi and Peters, 1998; Davis, 1998; Hull and Grefenstette, 1996). One of the important tasks of cross-language IR is to translate queries from one language to another. The original query and the translated query are then used to match documents in both the source and target languages. Target language documents are either glossed or translated by other systems. According to (Grefenstette, 1998), three main problems of query translations are:

1. generating translation candidates,
2. weighting translation candidates, and

3. pruning translation alternatives for document matching.

In cross-language IR, key word disambiguation is even more critical than in monolingual IR (Ballesteros and Croft, 1998) since the wrong translation can lead to a large amount of garbage documents in the target language, in addition to the garbage documents in the source language. Once again, we believe that sentence-based queries provide more information than mere key words in cross-language IR.

In both monolingual IR and cross-language IR, the query sentence or key words are assumed to be *consistently* in one language only. This makes sense in cases where the user is more likely to be a monolingual person who is looking for information in any language. It is also easier to implement a monolingual search engine. However, we suggest that the typical user of a cross-language IR system is likely to be bilingual to some extent. Most Web users in the world know some English. In fact, since English still constitutes 88% of the current web pages, speakers of another language would like to find English contents as well as contents in their own language. Likewise, English speakers might want to find information in another language. A typical example is a Chinese user looking for the information of an American movie, s/he might not know the Chinese name of that movie. His/her query for this movie is likely to be in **mixed language**.

Mixed language query is also prevalent in spoken language. We have observed this to be a common phenomenon among users of our SALSA system. The colloquial Hong Kong language is Cantonese with mixed English words. In general, a mixed language consists of a sentence mostly in the *primary language* with some words in a *secondary language*. We are interested in translating such mixed language queries into monolingual queries unambiguously.

In this paper, we propose a mixed language query disambiguation approach which makes use of the co-occurrence information of words between those in the primary language and those in the secondary language. We describe the overall methodology in Section 2. In Sections 2.1-3, we present the solutions to the three disambiguation problems. In Section 2.3 we present three different discriminative features

for disambiguation, ranging from the baseline model (Section 2.3.1), to the voting scheme (Section 2.3.2), and finally the 1-best model (Section 2.3.3). We describe our evaluation experiments in Section 3, and present the results in Section 4. We then conclude in Section 5.

2 Methodology

Mixed language query translation is halfway between query translation and query disambiguation in that not all words in the query need to be translated.

There are two ways to use the disambiguated mixed language queries. In one scenario, all secondary language words are translated unambiguously into the primary language, and the resulting monolingual query is processed by a general IR system. In another scenario, the primary language words are converted into secondary language and the query is passed to another IR system in the secondary language. Our methods allows for both general and cross-language IR from a mixed language query.

To draw a parallel to the three problems of query translation, we suggest that the three main problems of mixed language disambiguation are:

1. generating translation candidates in the primary language,
2. weighting translation candidates, and
3. pruning translation alternatives for query translation.

Co-occurrence information between neighboring words and words in the same sentence has been used in phrase extraction (Smadja, 1993; Fung and Wu, 1994), phrasal translation (Smadja et al., 1996; Kupiec, 1993; Wu, 1995; Dagan and Church, 1994), target word selection (Liu and Li, 1997; Tanaka and Iwasaki, 1996), domain word translation (Fung and Lo, 1998; Fung, 1998), sense disambiguation (Brown et al., 1991; Dagan et al., 1991; Dagan and Itai, 1994; Gale et al., 1992a; Gale et al., 1992b; Gale et al., 1992c; Shütze, 1992; Gale et al., 1993; Yarowsky, 1995), and even recently for query translation in cross-language IR as well (Ballesteros and Croft, 1998). Co-occurrence statistics is collected from either bilingual parallel and

non-parallel corpora (Smadja et al., 1996; Kupiec, 1993; Wu, 1995; Tanaka and Iwasaki, 1996; Fung and Lo, 1998), or monolingual corpora (Smadja, 1993; Fung and Wu, 1994; Liu and Li, 1997; Shütze, 1992; Yarowsky, 1995). As we noted in (Fung and Lo, 1998; Fung, 1998), parallel corpora are rare in most domains. We want to devise a method that uses only monolingual data in the primary language to train co-occurrence information.

2.1 Translation candidate generation

Without loss of generality, we suppose the mixed language sentence consists of the words $S = \{E_1, E_2, \dots, C, \dots, E_n\}$, where C is the only secondary language word¹. Since in our method we want to find the co-occurrence information between all E_i and C from a *monolingual* corpus, we need to translate the latter into the primary language word E_c . This corresponds to the first problem in query translation—translation candidate generation. We generate translation candidates of C via an online bilingual dictionary. All translations of secondary language word C , comprising of multiple senses, are taken together as a set $\{E_{c_i}\}$.

2.2 Translation candidate weighting

Problem two in query translation is to weight all translation candidates for C . In our method, the weights are based on co-occurrence information. The hypothesis is that the correct translations of C should co-occur frequently with the contextual words E_i and incorrect translation of C should co-occur rarely with the contextual words. Obviously, other information such as syntactical relationship between words or the part-of-speech tags could be used as weights too. However, it is difficult to parse and tag a mixed language sentence. The only information we can use to disambiguate C is the co-occurrence information between its translation candidates $\{E_{c_i}\}$ and E_1, E_2, \dots, E_n .

Mutual information is a good measure of the co-occurrence relationship between two words (Gale and Church, 1993). We first compute the mutual information between any word pair from a monolingual corpus in the primary language²

¹In actual experiments, each sentence can contain multiple secondary language words

²This corpus does not need to be in the same domain as the testing data

using the following formula, where E is a word and $f(E)$ is the frequency of word E .

$$MI(E_i, E_j) = \log \frac{f(E_i, E_j)}{f(E_i) * f(E_j)} \quad (1)$$

E_i and E_j can be either neighboring words or any two words in the sentence.

2.3 Translation candidate pruning

The last problem in query translation is selecting the target translation. In our approach, we need to choose a particular E_c from E_{c_i} . We call this pruning process **translation disambiguation**.

We present and compare three unsupervised statistical methods in this paper. The first baseline method is similar to (Dagan et al., 1991; Dagan and Itai, 1994; Ballesteros and Croft, 1998; Smadja et al., 1996), where we use the nearest neighboring word of the secondary language word C as feature for disambiguation. In the second method, we choose all contextual words as disambiguating feature. In the third method, the most discriminative contextual word is selected as feature.

2.3.1 Baseline: single neighboring word as disambiguating feature

The first disambiguating feature we present here is similar to the statistical feature in (Dagan et al., 1991; Smadja et al., 1996; Dagan and Itai, 1994; Ballesteros and Croft, 1998), namely the co-occurrence with neighboring words. We do not use any syntactic relationship as in (Dagan and Itai, 1994) because such relationship is not available for mixed-language sentences. The assumption here is that the most powerful word for disambiguating a word is the one next to it. Based on mutual information, the primary language target word for C is chosen from the set $\{E_{c_i}\}$. Suppose the nearest neighboring word for C in S is E_y , we select the target word E_{c_r} , such that the mutual information between E_{c_r} and E_y is maximum.

$$r = \operatorname{argmax}_i MI(E_{c_i}, E_y) \quad (2)$$

E_y is taken to be either the left or the right neighbor of our target word.

This idea is illustrated in Figure 1. MI1, represented by the solid line, is greater than MI2,

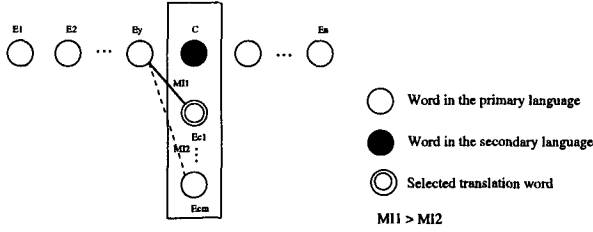


Figure 1: The neighboring word as disambiguating feature

represented by the dotted line. E_y is the neighboring word for C . Since MI_1 is greater than MI_2 , E_{c_1} is selected as the translation of C .

2.3.2 Voting: multiple contextual words as disambiguating feature

The baseline method uses only the neighboring word to disambiguate C . Is one or two neighboring word really sufficient for disambiguation?

The intuition for choosing the nearest neighboring word E_y as the disambiguating feature for C is based on the assumption that they are part of a phrase or collocation term, and that there is only one sense per collocation (Dagan and Itai, 1994; Yarowsky, 1993). However, in most cases where C is a single word, there might be some other words which are more useful for disambiguating C . In fact, such long-distance dependency occurs frequently in natural language (Rosenfeld, 1995; Huang et al., 1993).

Another reason against using single neighboring word comes from (Gale and Church, 1994) where it is argued that as many as 100,000 context words might be needed to have high disambiguation accuracy. (Shütze, 1992; Yarowsky, 1995) all use multiple context words as discriminating features. We have also demonstrated in our domain translation task that multiple context words are useful (Fung and Lo, 1998; Fung and McKeown, 1997).

Based on the above arguments, we enlarge the disambiguation window to be the entire sentence instead of only one word to the left or right. We use all the contextual words in the query sentence. Each contextual word “votes” by its mutual information with all translation candidates.

Suppose there are n primary language words in $S = E_1, E_2, \dots, C, \dots, E_n$, as shown in Figure 2, we compute mutual information scores

between all E_{c_i} and all E_j where E_{c_i} is one of the translation candidates for C and E_j is one of all n words in S . A mutual information score matrix is shown in Table 1. where MI_{jc_i} is the mutual information score between contextual word E_j and translation candidate E_{c_i} .

	E_{c_1}	E_{c_2}	...	E_{c_m}
E_1	$MI1c_1$	$MI1c_2$...	$MI1c_m$
E_2	$MI2c_1$	$MI2c_2$...	$MI2c_m$
...				
E_j	$MIjc_1$	$MIjc_2$...	$MIjc_m$
...				
E_n	$MInc_1$	$MInc_2$...	$MInc_m$

Table 1: Mutual information between all translation candidates and words in the sentence

For each row j in Table 1, the largest scoring MI_{jc_i} receives a vote. The rest of the row get zero's. At the end, we sum up all the one's in each column. The column i receiving the highest vote is chosen as the one representing the real translation.

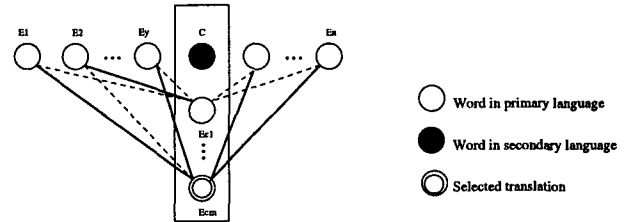


Figure 2: Voting for the best translation

To illustrate this idea, Table 2 shows that candidate 2 is the correct translation for C . There are four candidates of C and four contextual words to disambiguate C .

	E_{c_1}	E_{c_2}	E_{c_3}	E_{c_4}
E_1	0	1	0	0
E_2	1	0	0	0
E_3	0	0	0	1
E_4	0	1	0	0

Table 2: Candidate 2 is the correct translation

2.3.3 1-best contextual word as disambiguating feature

In the above voting scheme, a candidate receives either a one vote or a zero vote from *all contex-*

tual words equally no matter how these words are related to C . As an example, in the query “Please show me the latest dianying/movie of Jacky Chan”, *the* and *Jacky* are considered to be equally important. We believe however, that if the most powerful word is chosen for disambiguation, we can expect better performance. This is related to the concept of “trigger pairs” in (Rosenfeld, 1995) and Singular Value Decomposition in (Shütze, 1992).

In (Dagan and Itai, 1994), syntactic relationship is used to find the most powerful “trigger word”. Since syntactic relationship is unavailable in a mixed language sentence, we have to use other type of information. In this method, we want to choose the best trigger word among all contextual words. Referring again to Table 1, $MI_{j_c_i}$ is the mutual information score between contextual word E_j and translation candidate E_{c_i} .

We compute the *disambiguation contribution ratio* for each context word E_j . For each row j in Table 1, the largest MI score $MI_{j_c_f}$ and the second largest MI score $MI_{j_c_s}$ are chosen to yield the contribution for word E_j , which is the ratio between the two scores

$$\text{Contribution}(E_j, E_{c_i}) = \frac{MI_{j_c_f}}{MI_{j_c_s}} \quad (3)$$

If the ratio between $MI_{j_c_f}$ and $MI_{j_c_s}$ is close to one, we reason that E_j is not discriminative enough as a feature for disambiguating C . On the other hand, if the ratio between $MI_{i_e_f}$ and $MI_{i_e_s}$ is noticeably greater than one, we can use E_j as the feature to disambiguate $\{E_{c_i}\}$ with high confidence. We choose the word E_y with maximum contribution as the disambiguating feature, and select the target word E_{c_r} , whose mutual information score with E_y is the highest, as the translation for C .

$$r = \arg \max_i MI(E_y, E_{c_i}) \quad (4)$$

This method is illustrated in Figure 3. Since E_2 is the contextual word with highest contribution score, the candidate E_i is chosen that the mutual information between E_2 and E_{c_i} is the largest.

3 Evaluation experiments

The mutual information between co-occurring words and its contribution weight is ob-

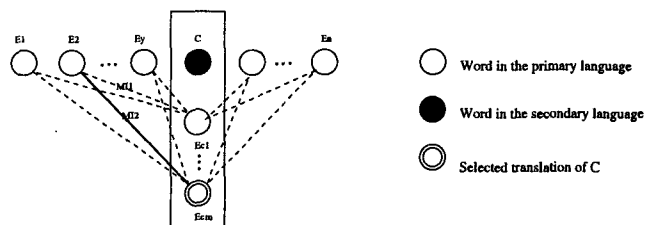


Figure 3: The best contextual word as disambiguating feature

tained from a monolingual training corpus—Wall Street Journal from 1987-1992. The training corpus size is about 590MB. We evaluate our methods for mixed language query disambiguation on an automatically generated mixed-language test set. No bilingual corpus, parallel or comparable, is needed for training.

To evaluate our method, a mixed-language sentence set is generated from the monolingual ATIS corpus. The primary language is English and the secondary language is chosen to be Chinese. Some English words in the original sentences are selected randomly and translated into Chinese words manually to produce the testing data. These are the mixed language sentences. 500 testing sentences are extracted from the ARPA ATIS corpus. The ratio of Chinese words in the sentences varies from 10% to 65%.

We carry out three sets of experiments using the three different features we have presented in this paper. In each experiment, the percentage of primary language words in the sentence is incrementally increased at 5% steps, from 35% to 90%. We note the accuracy of unambiguous translation at each step. Note that at the 35% stage, the primary language is in fact Chinese.

4 Evaluation results

One advantage of using the artificially generated mixed-language test set is that it becomes very easy to evaluate the performance of the disambiguation/translation algorithm. We just need to compare the translation output with the original ATIS sentences.

The experimental results are shown in Figure 4. The horizontal axis represents the percentage of English words in the testing data and the vertical axis represents the translation accuracy. Translation accuracy is the ratio of the number of secondary language (Chinese) words disambiguated correctly over the number of all

secondary language (Chinese) words present in the testing sentences. The three different curves represent the accuracies obtained from the baseline feature, the voting model, and the 1-best model.

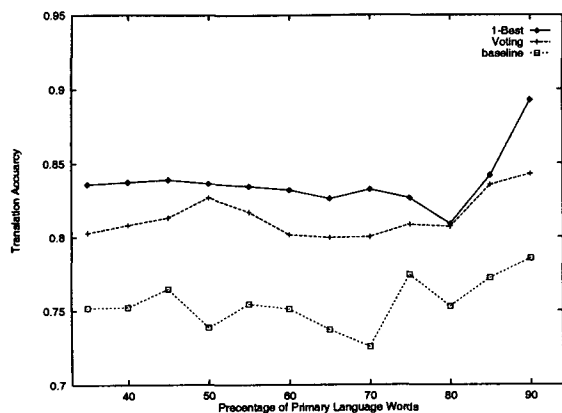


Figure 4: 1-best is the most discriminating feature

We can see that both voting contextual words and the 1-best contextual words are more powerful discriminant than the baseline neighboring word. The 1-best feature is most effective for disambiguating secondary language words in a mixed-language sentence.

5 Conclusion and Discussion

Mixed-language query occurs very often in both spoken and written form, especially in Asia. Such queries are usually in complete sentences instead of concatenated word strings because they are closer to the spoken language and more natural for user. A mixed-language sentence consists of words mostly in a primary language and some in a secondary language. However, even though mixed-languages are in sentence form, they are difficult to parse and tag because those secondary language words introduce an ambiguity factor. To understand a query can mean finding the matched document, in the case of Web search, or finding the corresponding semantic classes, in the case of an interactive system. In order to understand a mixed-language query, we need to translate the secondary language words into primary language *unambiguously*.

In this paper, we present an approach of mixed-language query disambiguation by using co-occurrence information obtained from a

monolingual corpus. Two new types of disambiguation features are introduced, namely voting contextual words and 1-best contextual word. These two features are compared to the baseline feature of a single neighboring word. Assuming the primary language is English and the secondary language Chinese, our experiments on English-Chinese mixed language show that the average translation accuracy for the baseline is 75.50%, for the voting model is 81.37% and for the 1-best model, 83.72%.

The baseline method uses only the neighboring word to disambiguate C . The assumption is that the neighboring word is the most semantic relevant. This method leaves out an important feature of nature language: long distance dependency. Experimental results show that it is not sufficient to use only the nearest neighboring word for disambiguation.

The performance of the voting method is better than the baseline because more contextual words are used. The results are consistent with the idea in (Gale and Church, 1994; Shütze, 1992; Yarowsky, 1995).

In our experiments, it is found that 1-best contextual word is even better than multiple contextual words. This seemingly counter-intuitive result leads us to believe that choosing the most discriminative single word is even more powerful than using multiple contextual word equally. We believe that this is consistent with the idea of using “trigger pairs” in (Rosenfeld, 1995) and Singular Value Decomposition in (Shütze, 1992).

We can conclude that sometimes long-distance contextual words are more discriminant than immediate neighboring words, and that multiple contextual words can contribute to better disambiguation. Our results support our belief that natural sentence-based queries are less ambiguous than keyword based queries. Our method using multiple disambiguating contextual words can take advantage of syntactic information even when parsing or tagging is not possible, such as in the case of mixed-language queries.

Other advantages of our approach include: (1) the training is unsupervised and no domain-dependent data is necessary, (2) neither bilingual corpora or mixed-language corpora is needed for training, and (3) it can generate

monolingual queries in both primary and secondary languages, enabling true cross-language IR.

In our future work, we plan to analyze the various “discriminating words” contained in a mixed language or monolingual query to find out which class of words contribute more to the final disambiguation. We also want to test the significance of the co-occurrence information of all contextual words *between themselves* in the disambiguation task. Finally, we plan to develop a general mixed-language and cross-language understanding framework for both document retrieval and interactive tasks.

References

- AskJeeves. 1998. <http://www.askjeeves.com>.
- Lisa Ballesteros and W. Bruce Croft. 1998. Resolving ambiguity for cross-language retrieval. In *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 64–71, Melbourne, Australia, August.
- P. Brown, J. Lai, and R. Mercer. 1991. Aligning sentences in parallel corpora. In *Proceedings of the 29th Annual Conference of the Association for Computational Linguistics*.
- Ido Dagan and Kenneth W. Church. 1994. Termight: Identifying and translating technical terminology. In *Proceedings of the 4th Conference on Applied Natural Language Processing*, pages 34–40, Stuttgart, Germany, October.
- Ido Dagan and Alon Itai. 1994. Word sense disambiguation using a second language monolingual corpus. In *Computational Linguistics*, pages 564–596.
- Ido Dagan, Alon Itai, and Ulrike Schwall. 1991. Two languages are more informative than one. In *Proceedings of the 29th Annual Conference of the Association for Computational Linguistics*, pages 130–137, Berkeley, California.
- M. Davis. 1998. Free resources and advanced alignment for cross-language text retrieval. In *Proceedings of the 6th Text Retrieval Conference (TREC-6)*, NIST, Gaithersburg, MD, November.
- Laura DiDio. 1997. Os/2 let users talk back to 'net. page 12.
- ElectricMonk. 1998. <http://www.electricmonk.com>.
- Pascale Fung and Yuen Yee Lo. 1998. An IR approach for translating new words from non-parallel, comparable texts. In *Proceedings of the 36th Annual Conference of the Association for Computational Linguistics*, pages 414–420, Montreal, Canada, August.
- Pascale Fung and Kathleen McKeown. 1997. Finding terminology translations from non-parallel corpora. In *The 5th Annual Workshop on Very Large Corpora*, pages 192–202, Hong Kong, Aug.
- Pascale Fung and Dekai Wu. 1994. Statistical augmentation of a Chinese machine-readable dictionary. In *Proceedings of the Second Annual Workshop on Very Large Corpora*, pages 69–85, Kyoto, Japan, June.
- Pascale Fung, CHEUNG Chi Shuen, LAM Kwok Leung, LIU Wai Kat, and LO Yuen Yee. 1998a. A speech assisted online search agent (salsa). In *ICSLP*.
- Pascale Fung, CHEUNG Chi Shuen, LAM Kwok Leung, LIU Wai Kat, LO Yuen Yee, and MA Chi Yuen. 1998b. SALSA, a multilingual speech-based web browser. In *The First AEARU Web Technology Workshop*, Nov.
- Pascale Fung. 1998. A statistical view of bilingual lexicon extraction: from parallel corpora to non-parallel corpora. In *Proceedings of the Third Conference of the Association for Machine Translation in the Americas*, Pennsylvania, October.
- William A. Gale and Kenneth W. Church. 1993. A program for aligning sentences in bilingual corpora. *Computational Linguistics*, 19(1):75–102.
- William A. Gale and Kenneth W. Church. 1994. Discrimination decisions in 100,000 dimensional spaces. *Current Issues in Computational Linguistics: In honour of Don Walker*, pages 429–550.
- W. Gale, K. Church, and D. Yarowsky. 1992a. Estimating upper and lower bounds on the performance of word-sense disambiguation programs. In *Proceedings of the 30th Conference of the Association for Computational Linguistics*. Association for Computational Linguistics.
- W. Gale, K. Church, and D. Yarowsky. 1992b.

- Using bilingual materials to develop word sense disambiguation methods. In *Proceedings of TMI 92*.
- W. Gale, K. Church, and D. Yarowsky. 1992c. Work on statistical methods for word sense disambiguation. In *Proceedings of AAAI 92*.
- W. Gale, K. Church, and D. Yarowsky. 1993. A method for disambiguating word senses in a large corpus. In *Computers and Humanities*, volume 26, pages 415–439.
- Gregory Grefenstette, editor. 1998. *Cross-language Information Retrieval*. Kluwer Academic Publishers.
- Xuedong Huang, Fileno Alleva, Hisao-Wuen Hong, Mei-Yuh Hwang, Kai-Fu Lee, and Ronald Rosenfeld. 1993. The SPHINX-II speech recognition system: an overview. *Computer, Speech and Language*, pages 137–148.
- David A. Hull and Gregory Grefenstette. 1996. A dictionary-based approach to multilingual information retrieval. In *Proceedings of the 19th International Conference on Research and Development in Information Retrieval*, pages 49–57.
- Julian Kupiec. 1993. An algorithm for finding noun phrase correspondences in bilingual corpora. In *Proceedings of the 31st Annual Conference of the Association for Computational Linguistics*, pages 17–22, Columbus, Ohio, June.
- Xiaohu Liu and Sheng Li. 1997. Statistic-based target word selection in English-Chinese machine translation. *Journal of Harbin Institute of Technology*, May.
- Chi Yuen Ma and Pascale Fung. 1998. Using English phoneme models for Chinese speech recognition. In *International Symposium on Chinese Spoken language processing*.
- D.W. Oard. 1997. Alternative approaches for cross-language text retrieval. In *AAAI Symposium on cross-language text and speech retrieval*. American Association for Artificial Intelligence, Mar.
- Eugenio Picchi and Carol Peters. 1998. Cross-language information retrieval: a system for comparable corpus querying. In Gregory Grefenstette, editor, *Cross-language Information Retrieval*, pages 81–92. Kluwer Academic Publishers.
- Lau Raymond. 1997. Webgalaxy : Beyond point and click - a conversational interface to a browser. In *Computer Networks & ISDN Systems*, pages 1385–1393.
- Rony Rosenfeld. 1995. *A Corpus-Based Approach to Language Learning*. Ph.D. thesis, Carnegie Mellon University.
- Hinrich Shütze. 1992. Dimensions of meaning. In *Proceedings of Supercomputing '92*.
- Frank Smadja, Kathleen McKeown, and Vasileios Hatzivassiloglou. 1996. Translating collocations for bilingual lexicons: A statistical approach. *Computational Linguistics*, 21(4):1–38.
- Frank Smadja. 1993. Retrieving collocations from text: Xtract. *Computational Linguistics*, 19(1):143–177.
- Kumiko Tanaka and Hideya Iwasaki. 1996. Extraction of lexical translations from non-aligned corpora. In *Proceedings of COLING 96*, Copenhagen, Denmark, July.
- Dekai Wu. 1995. Grammarless extraction of phrasal translation examples from parallel texts. In *Proceedings of TMI 95*, Leuven, Belgium, July. Submitted.
- D. Yarowsky. 1993. One sense per collocation. In *Proceedings of ARPA Human Language Technology Workshop*, Princeton.
- D. Yarowsky. 1995. Unsupervised word sense disambiguation rivaling supervised methods. In *Proceedings of the 33rd Conference of the Association for Computational Linguistics*, pages 189–196. Association for Computational Linguistics.
- Victor Zue. 1995. Spoken language interfaces to computers: Achievements and challenges. In *The 33rd Annual Meeting of the Association of Computational Linguistics*, Boston, June.