

# **eftab200**

## **Using Linguistic Tools and Resources in Cross-Language Retrieval**

Carol Peters\* and Eugenio Picchi\*\*

\*Istituto di Elaborazione della Informazione - Consiglio Nazionale delle Ricerche  
Via S. Maria, 46, 56126 PISA, Italy  
carol@iei.pi.cnr.it

\*\*Istituto di Linguistica Computazionale - Consiglio Nazionale delle Ricerche  
Via della Faggiola, 32, 56126 PISA, Italy  
picchi@ilc.pi.cnr.it

### **Abstract**

A system to process bilingual/multilingual text corpora is described. The system includes components for cross-language querying on parallel (i.e. translation equivalent) and comparable (i.e. domain-specific) collections of texts in more than one language. Both sets of procedures are dependent on lexical resources (bilingual lexical databases) and linguistic tools (morphological procedures). The system was originally designed to meet the requirements of various types of contrastive language studies. However, we are now studying applications to cross-language retrieval.

### **Background**

In the last few years, natural language processing (NLP) techniques and tools have been incorporated into information retrieval (IR) systems with varying degrees of success (Smeaton, 1992). The recent emergence of the field of Cross-Language Information Retrieval as an independent area of interest has clearly reinforced this trend. In order to be successful, cross-language applications frequently need access to methodologies and resources that were originally studied and constructed for NLP purposes, such as morphological analysers and generators, computational lexicons, various kinds of procedures for text analysis, etc. The integration of such resources into typical IR processes implies an exchange of know-how and viewpoints between the two disciplines.

### **The PiSystem**

At the Institute for Computational Linguistics (ILC-CNR) in Pisa, there has been intensive work over the last decade on the development of an integrated complex set of mono- and bilingual lexicon and text management and analysis tools, known as the PiSystem. These tools are designed to meet the needs of all kinds of literary and linguistic text processing tasks. The core component of the system is the DBT (Textual Database) search engine. The DBT system has been implemented in various versions to process and analyse different kinds of structured and unstructured texts. There is also a client-server version running on Internet, known as DBTNET. A morphological engine and a part-of-speech tagger and lemmatizer for Italian are also provided. (For full details and demos, see our Web site: <http://www.ilc.pi.cnr.it/dbt/episystem/index.htm>).

Part of the PiSystem consists of a set of components designed for multilingual or cross-language applications. This includes a system for the acquisition, management and querying of mono- and bilingual lexical databases (LDBs), morphological data and rules for Italian, English, French, Latin, and a system for bilingual/multilingual text management, which has separate procedures for cross-language querying on both Parallel and Comparable Text Corpora. Both sets of procedures use other components of the PiSystem:

- DBT
- Bilingual Lexical Database
- Morphological Analysers and Generators

We define *ParallelCorpora* as collections of translationally equivalent texts regarding general language or sub-languages of different types (e.g. for a given author, a particular domain, etc.), and *Comparable Corpora* as homogeneous sets of texts from pairs or multiples of languages with the same communicative function; they must share certain basic features, which could be period, author, style, genre, register, but generally they refer to the same domain; they thus regard sub-languages rather than general languages. The two types of corpora provide different kinds of contrastive data: parallel corpora provide data on translation equivalents; comparable corpora give information on natural language lexical equivalents within a given domain<sup>1</sup>.

Our corpus procedures were initially developed to process Italian/English texts but are extendible to other languages. We are now adding modules to handle French texts. The system was originally studied with a number of human-oriented applications in mind: bilingual lexicography; language learning activities; translating and translation studies; cross-linguistic studies. We are now studying possible applications to CLIR activities.

## Parallel Text System

Our parallel text system has been described elsewhere (Marinai et al, 1991, 1992) and we will not go into details concerning its implementation here. Its distinguishing feature is that it is based on the use of linguistic resources rather than the statistical-distributional data employed by the most well known procedures for this type of application (see, for example, Warwick and Russell, 1990, Gale and Church, 1993, Brown et al, 1993). Thus, it does not perform a sentence-based alignment but uses a bilingual electronic dictionary and morphological analysers and generators for Italian and English to link texts and construct parallel contexts through the recognition of translation equivalents in pairs of texts.

Our procedures function in two distinct stages:

1. Links are created between translation equivalent pairs in texts being processed. This stage is performed just once for each new pair of parallel text added to the archives;
2. The query system uses the links to construct the parallel contexts, in real time, for any form or co-occurrence of forms contained in the texts.

Full details on how the text linking algorithm operates, how the 'search zone' (area of L2 corpus searched for L2 equivalents of a given L1 word) is calculated, how false links are eliminated, and under what conditions failure is signalled, are given in Marinai et al. (1992). It is sufficient to note that, from the user viewpoint, this operation is simple, rapid and, once a few preliminary instructions have been given, automatic. There is no need to pre-process the texts, manually indicating, for example, matching sections of the text. Neither is it important that every possible translation match between two texts is recognised; the translation links are just the means to construct the parallel contexts as accurately as possible. However, the fact that the system uses and evidences information on translations derived from a general purpose bilingual dictionary facilitates the recognition and retrieval of 'new' information, i.e. real world translations for words or expressions not given in the bilingual dictionary (as shown in Figure 1 below).

When the system is queried to retrieve parallel contexts for any word or combination of words from the corpora, the L1 context is constructed with the searched word(s) in the centre; these words are highlighted and during creation of the context, any translation links to the L2 text associated with the other words in the context are read. If the searched words have an associated link, this is used to identify directly the corresponding word(s) in

QUERY = sicurezza 5 Contexts Found  
 {I} Laprese come un avvertimento e, voltatosi indietro verso la casa nella quale era entrata la giovane, per avere la sicurezza di non essere osservato, traverso' dicorsa la strada senza aspettare. L'ansia e l'andatura veloce lo facevano ansimare - =FE - I-Dublin6.39  
 {E} light rain fell. He took them as a warning and glancing back towards the house which the young woman had entered to see that he was not observed, he ran eagerly across the road. Anxiety and his swift run made him pant - =FE - E-Dublin6.105

<sup>1</sup> Unfortunately, there is still no general consensus between the applied linguistics and the computational linguistics communities on the definition of these two terms. The applied linguists tend to use them differently from us. Hartmann

{I} Gabriel risecon un certo nervosismo e dovette dare un tocco al cravattino peracquistare **sicurezza**, mentre zia Kate si piegava in due dalle risa tanto le era piaciuto loscherzo =FE - I-Dublin7.128  
 {E} Gabriel **laughed** nervously and patted histie reassuringly **while** Aunt Kate nearly doubledherself, so heartily did she enjoy the joke. =FE - E-Dublin7.152  
 {I} Bessie che si occupa di loro." "Benone", ripeté zia Kate. "E' unabella **sicurezza**, una ragazza come quella, una sucui poter contare! C'è quella Lily, invece, =FE - I-Dublin15.232  
 {E} "Besides, Bessie will look after them." "To be sure", said Aunt**Kate** again. "What a comfort it is to have **agirl** like that, one you can depend on! There'sthat Lily, =FE - E-Dublin15.262  
 - DBT -- E.Picchi-----

**Figure 1: ParallelContexts for Italian term *sicurezza* in "TheDubliners"**

the L2 text, which will also behighlighted and used as the central point for the construction of the L2 context; other words that have been linked in thepaired contexts can be optionally evidenced in a different colour. The sizeof the context and the maximum distance between cooccurrences of items searched can be specified by the user. Figure 1 shows the firstthree parallel contexts found for the Italian word *sicurezza* in a parallel corpus consisting of the original text of "The Dubliners" byJames Joyce plus an Italian translation.

When there is no directly linked L2 formfor the L1 word being searched (as is the case in the figure), then all the linksfor words found in the source context are used to calculate an 'average'value which identifies the central point around which the relative L2context will be constructed. The calculation of this 'average' value allows for the possibility of unevenconcentrations of matched words in the contexts. 'Wrong' links betweenpotentially falsely recognised translation equivalents which disturbcontext calculation are identified and eliminatedby the query system. The two linked forms which are closest to the pointcalculated as the middle of the target text context are evidenced in adifferent colour, as indicators of the likely position of the translationequivalent. In the figure, the word searched and the indicators of its position in the target context are shownin bold. Thus, for the first parallel context, we see that the position ofthe L2 equivalent of *sicurezza* ("*per avere la sicurezza*" is translated by "to see that") is found between "young" and"observed"; in the second, (where the concept of *sicurezza*is translated by the adverb "reassuringly") the translation is indicated aslying between "laughed" and "while"; in the third (*sicurezza* = comfort) it falls between "Kate" and "girl". In the first pair ofcontexts only, each word for which a translation link to a word in theother language text has been created is underlined in order to show theinformation the system uses to construct the parallel context and to indicate the position of the target languageequivalent.

The system has been tested on all kindsof texts: scientific articles, extracts from text books, on-flightmagazines, short stories, novels, poetry.It is easy to evaluate the results; the system is successful if, for eachoccurrence of any form or occurrence of forms in the texts in one language,correct parallel contexts are constructed for all the translation equivalents contained in the texts in thesecond language. The only factor that may affect performance is if the userrequests a very small context size. The system default value is 25 tokenseach side of the keyword(s) (the items searched); it is advisable not to request contexts smaller than 15tokens as, in this case, the system may not have sufficient links tocalculate the parallel contexts accurately.

It is now our intention to testthis system in a translation-training application.

The approach described has beencriticised as being costly in terms of resources and only extendible toother pairs of languages if the necessary lexical/linguistic data areavailable. Our answer is that, in the first place, such resources(dictionaries and morphologies) are now widely available in thecomputational linguistics community for most European languages and, ingrowing numbers, for non-European languages. It wouldbe foolish not to use them. Secondly, it is false to believe that alignmentprocedures can be considered as totally language-independent. McEnery andOakes (1996), for instance, show how alignment methods of the typeexemplified by Gale and Church (*op.cit.*) are both language and domain dependent. These authors report recentliterature in which various suggestions are made as to how this type ofalignment can be improved by introducing some surface linguistic knowledgesuch as the notion of language specific cognate words, ie pairs of tokens in a given language which share'obvious' phonologic or orthographic and semantic properties. They go on to describe methods to employ this kind of cognate informationin order to enhance statistical sentence alignment results, although theyalso note that the incorporation of too much cognate information willresult in noise and a degradation of results.

---

(1995), forinstance, has proposed a distinction between bi-texts, for translationallylinked texts, and parallel texts for texts that are functionally similar insituational motivation and rhetorical structure.

The recent trend in parallel corpus processing is thus to move towards the use of other than purely statistical data; this appears to support our claim that if lexical and morphological components are used in procedures for bilingual text alignment, performance in terms of retrieval precision can be improved.

## Comparable Text System

We recently decided to extend the scope of our bilingual corpus system by including a set of procedures for the analysis and extraction of significant data from comparable text archives. The aim was different but our interests were still oriented towards human-oriented language learning activities: while with our parallel system the user can retrieve examples of specific instances of how a given word or expression has been translated in another language, depending on context, argument, stylistic considerations, etc., using the comparable system, he/she will also be able to look for natural language examples of L2 lexical equivalents of a given word or expression in L1, independently of any direct translation link. By definition, comparable text archives regard special domains or sublanguages; they are thus of particular interest for studies or applications regarding terminology and for technical rather than literary texts.

Our procedures operate on sets of comparable texts in two different languages. We are currently working on Italian and English texts, and so far all work has been focused on nouns; in sub-language texts, it is mainly the nouns that bear the weight of topic-specificity, i.e. the technical message. The approach is based on the assumptions that (i) words acquire sense from their context, (ii) words used in a similar way throughout a sub-language or special domain corpus will be semantically similar. It follows that, if it is possible to establish equivalences between several items contained in two different contexts, there is a high probability that the two contexts themselves are to some extent similar. It is important to stress that our aim is not to retrieve precise equivalences in L2 of the L1 term under examination, but to isolate the set of contexts in the L2 corpora that has the highest probability of providing L2 correspondences to the L1 input. Given a particular term or set of terms found in the texts in one language (L1), we attempt to identify contexts which treat the same argument in the texts of the second language (L2). To do this, we isolate the vocabulary related to that term in the L1 corpus - hypothesising that the word will be surrounded by a similar vocabulary in L2.

A term, T, is thus selected in the one set of texts (denominated as L1 - either set can be chosen as L1). For each occurrence of T in the L1 set of texts, the system constructs a context window containing T plus up to 'n' lexically significant words appearing to the right and left of T, but within the same phrase, i.e. strong punctuation marks (full stops and semi-colons) act as break points in the construction of these contexts. The value for 'n' is set by the user. Words contained in a stop list are not counted. This list includes functional words such as articles, pronouns, prepositions, and also highly frequent, insignificant words which would create noise. The stop list can be modified by the user so that certain frequent terms specific to the particular domain can be eliminated if necessary to improve performance. In the current version of the system, we are accepting just nouns and verbs as being relevant for our purposes; we have made tests in which we accepted only nouns but we appeared to lose some significant information. We have also tested with varying values for 'n' in order to try to establish experimentally the optimum size of the context window in terms of significance of results weighed against processing times (clearly, the larger the context window, the longer the time needed to calculate the significant vocabulary for a given term).

For each co-occurrence of our keyword T in the context windows, morphological procedures identify the source lemma. The set of significant words found in the context windows for T make up the vocabulary,  $V_1$ , that is considered to characterise T in the particular L1 corpus being analysed. The frequencies of the co-occurrences of T are then computed and to each element of  $V_1$  is assigned its mutual information value which measures the significance of the correlation between the  $V_1$  item and T, i.e. the relative frequency of the  $V_1$  item as a collocate of T is measured against its overall frequency in the corpus in order to identify how strongly it is related to T (see Church and Hanks 1990). Using the MI index as an ordering element, we list  $V_1$  in order of decreasing significance and set a threshold below which terms in  $V_1$  are not considered relevant and can be ignored. Figure 2 shows the significant collocates for the Italian noun *libertà* found in a set of comparable English and Italian parliamentary debates. There were 198 occurrences of *libertà*.

0000000	16	
500.000	198	LIBERTA' (freedom)
11.259	4	CONDIZIONARE (to condition)
10.094	30	FONDAMENTALE (fundamental)
9.358	5	ESPRESSIONE (expression)
8.965	5	SINDACARE (to inspect/control)
8.722	3	EFFETTIVO (effective)
8.619	4	PIENA PIENARE (full to fill)
8.573	4	INDIVIDUARE (to identify)
8.550	4	BENEFICIARE (to benefit by/from)
8.204	3	LIMITAZIONE (limitation)
7.696	3	GARANZIA (guarantee)
6.672	5	PRINCIPIO PRINCIPIARE (principle to begin)
6.155	4	CITTADINO (citizen)
6.040	5	RISPETTO RISPETTARE (respect to respect)
5.975	6	POLITICA (politics/policy)
5.746	5	TRATTARE TRATTATO (to treat treaty)

**Figure 2: libertà - 198 occurrences - 16 significant collocates**

For each collocate, the first column shows the MI value, the second its frequency value, ie the number of times the collocate was found in the context windows for *libertà*; in this case, we accepted both nouns and verbs and 'n' was set at 5. An indication of the meaning of the collocates in English has been given between brackets for easier understanding.

Next, using our lexical tools (e.g. morphological analysers and generators and bilingual lexical database), we construct an equivalent vocabulary (V<sub>2</sub>) in L2 of translation equivalents for the L1 set of cooccurrences (V<sub>1</sub>), ie for each element of V<sub>1</sub> we create a set of L2 translation equivalents, denoted as L2 translation blocks.

Each block contains the entire set of translations supplied by the bilingual lexical database for a member of the L1 vocabulary (no distinction is made for sense), together with all possible forms for each translation (generated by the morphological procedure). For example the L2 translation block for the Italian lemma *garanzia* includes the English forms *guarantee, guarantees, security, securities, surety, sureties*. To each translation block, we assign a value equal to the MI Index of the L1 term represented by this translation block. These values are used to assign weights to the translation blocks to represent the probability of occurrence in the L2 texts of any of the members of that particular translation block when searching for expressions regarding our keyword, T. Direct translations of the term itself are also assigned an arbitrarily high value as being the most probable L2 representative of T. An L2 stoplist is also applied at this point, again in order to eliminate as much noise as possible from the items contained in the translation blocks; basically, we eliminate very common L2 words.

The procedure then searches the L2 corpus in order to identify words or expressions that can be considered in some way lexically equivalent to our selected term in the L1 texts. This is done by searching for those contexts in L2 in which there is a significant presence of the L2 vocabulary for T. The significance is determined on the basis of a statistical procedure; this procedure uses the number of V<sub>2</sub> items found in the context and the weights assigned to them in order to assess the probability that any given L2 cooccurrence represents a lexically equivalent context for T, and to establish thresholds of acceptability.

Comparable Contexts		
6	535.181930	1) is also a result of the fact that international *rules* requiring strict <b>safety</b> standards for *passenger* vessels apply only to those operating internationally. Because *FXAC93207ENC.0003.01.00".30
5	530.6161000	2) Council Directive of 30 November on the minimum *health* and <b>safety</b> requirements for the use by *workers* of *personal* protective equipment at the work-place (1) * "FXAC93207ENC.0042.01.00".22

5	528.9811000	14)	27 October 1992) \\ (93/C 327/22)\Q\	*Subject*: Proposal for a *Council* *regulation* on <b>security</b> *measures* applicable to classified information produced or * "FXAC93327ENC.0013.02.00".9
5	528.1311000	31)	purpose. They have given their full backing to *Resolution* 787 of the United *Nations* <b>Security</b> *Council*	which *stepped* up sanctions against the Federal Republic of* "FXAC93350ENC.0040.03.00".25
5	37.766 223	45)	Communities (12 March 1993) \\(93/C 280/85) \Q\	*Subject*: *Council* *Directive* laying down *health**rules* for the production and placing on the market of raw milk, heat* "FXAC93280ENC.0043.03.00".9
5	37.432 217	46)	to an installation using ionizing radiation. To *increase* the *health* *protection* of such *workers* the *Council* adopted, on 4 December 1990, Directive 90/641/	* "FXAC93065ENC.0017.01.00".75
5	34.725216	47)	States and the IAEA in the Commission Standing *Working* Group on the *transport* of *nuclear**materials* cover the *need* for mutual information between the parties. Since	* "FXAC93327ENC.0024.01.00".26
5	34.132 141	48)	for a Directive on the reorganization of working *time* setting minimum *requirements* at Community *level* for *protecting* the *health* and safety of workers. The provisions	* "FXAC93086ENC.0036.01.00".38
4	524.1361000	52)	the G7 summit in Munich that the design of RBMK *nuclear* *reactors* does not offer the same <b>safety</b> *guarantees* as do nuclear reactors of a more modern design. The	* "FXAC93327ENC.0028.02.00".28
4	30.083 214	205)	using ionizing radiation. To increase the *health* *protection* of such *workers* the *Council* adopted, on 4 December 1990, Directive 90/641/Euratom on	* "FXAC93065ENC.0017.01.00".75
----- DBT --E.Picchi -----				
The comparable contexts are ordered by (i) no. of significant collocates, (ii) presence of direct translations of the term searched, (iii) MI value, (iv) sum of frequency values. Column 4 shows their ranking order.				

**Figure 3: Comparable contexts for *sicurezza* in Italian/English parliamentary texts**

Although it is clear that the process of translating the L1 vocabulary for T into L2 introduces a number of irrelevant terms (all dictionary provided translations are accepted) and only some of these are eliminated by the L2 stop list, this does not normally affect the results as, if an L2 context is to be accepted as representative of a given L1 term, it is necessary for a number of items from the L2 vocabulary for T to be present.

The results are written in a file and ranked in descending order according to (i) the number of items in the context coming from the L2 vocabulary, (ii) whether a direct translation of the term being searched is included, (iii) the sum of the MI values associated with these items. The file of results can be displayed on the screen, saved, or printed out for further consultation. The user can also enlarge a selected context by clicking on it so that he/she can refer to the entire piece of text to which it belongs in the underlying corpus.

Figure 3 shows the results of a query on our comparable corpus for the Italian noun *sicurezza*. For reasons of space, we have printed out only 10 contexts, just to give an idea of the kind of results we obtain using this method. The first context contained 6 items from the L2 vocabulary for the term being searched including a direct translation of it; the next 43 contexts (2-44) contained 5 items from the L2 vocabulary, including direct translations - this was generally "safety" but in contexts, 14 and 31 we find "security"; Nos. 45-48 show examples of contexts in which there was no direct translation of the term itself, just the presence of 5 items from the L2 vocabulary. The reader can judge for him/herself to what extent he/she feels that the context represents the concepts of safety/security in this corpus. Context 52 was the first context containing just 4 items from the L2 vocabulary including "safety" and No. 245 was the first context with four L2 vocabulary terms which did not include a direct translation of *sicurezza*.

Our test corpus up until now has been a set of parliamentary debates in English and Italian of approximately 1 million word forms per language. This corpus has been useful for testing but is not entirely satisfactory for our purposes; although stylistically homogenous, the lexicon is rather too general and thus not suitable for studies on terminology as was our original intention; first results of testing on this corpus were given in Picchi, Peters (1996). We have now made some trial runs on a set of newspaper articles from the same year: *Corriere della Sera* (approx. 300,000 word forms) and the *Independent* (approx. 600,000 word forms). As expected, we find that in

order to obtain significant results the texts must be either highly homogeneous, or of large dimensions. In the case of the newspaper articles, the corpus was not sufficiently large to give us many interesting terms with a reasonably high frequency in order to calculate a significant vocabulary. (In fact, this collection of texts had been collected for a different purpose: to study the behaviour of certain neologisms over languages). In the next two figures, we give the results obtained for the Italian term *lavoro* (translated by the dictionary as: work, job, task, labour), one of the most frequent terms found in these newspaper archives.

DBT - Comparable texts		from
<b>Corriere della Sera 1994</b>		
163	lavoro	
0000000	16	
500.000	163	LAVORO (work, labour, job, task)
11.772	4	SUBORDINATO (subordinate)
10.100	5	AUTONOMO (autonomous)
8.907	3	MASTELLA (Mastella - Italian Minister of Work)
8.262	8	DIPENDENTE (employee)
8.203	3	CONTRATTO (contract)
7.794	3	RELAZIONE (report)
7.734	6	CAPITALE (capital)
7.681	3	OCCUPAZIONE (occupation)
7.563	8	MINISTRO (ministry)
7.526	3	FINANZA (finance)
7.478	3	ESIGENZA (necessity)
7.302	3	ORGANIZZAZIONE (organization)
7.145	3	PRINCIPALE (principal)
7.109	3	IMPRESA (agency)
6.940	3	SINDACATO (union)
6.938	6	TRATTA (trade)
6.879	9	POSTO (position)

**Figure 4: *lavoro* - 163 occurrences - 18 significant collocates**

DBT - Comparable texts		from
<b>The Independent 1994</b>		
4	521.295 176	1) in the home. Inability to switch off from <b>*work*</b> , <b>*experienced*</b> by 83 per cent, may be a <b>*major*</b> stress on MPs' family <b>*relationships*</b> . Looking for an <b>*IND Health:</b>
4	521.035 176	2) team-working ability, motivation and drive. We <b>*need*</b> graduates to <b>*work*</b> in the following broad skill areas: <b>*business*</b> analysts, engineers, <b>*finance*</b> <b>*IND Enterprise 94:</b>
4	520.757 185	3) told him, effectively been defeated for the <b>*Labour*</b> leadership, as a result of machinations by a group of <b>*trade*</b> <b>*union*</b> barons who <b>*controlled*</b> the <b>*IND Profile:</b>
4	27.823 19	4) designers and marketers versed in the latest <b>*business*</b> <b>*school*</b> techniques. <b>*Firms*</b> of the market seek out many <b>*employees*</b> with out-of-the-ordinary views <b>*INS</b>
4	27.435 15	5) than purely reactive. We are beginning to see <b>*trade*</b> <b>*unions*</b> as an important voice on the <b>*business*</b> pages, the personal <b>*finance*</b> pages, the health page <b>*IND</b>
4	27.387 15	6) Last month Phase closed after five <b>*issues*</b> . It faced a familiar small <b>*business*</b> problem - not enough <b>*capital*</b> to <b>*finance*</b> its early losses, <b>*IND</b>
4	27.387 15	7) after five issues. It faced a familiar small <b>*business*</b> <b>*problem*</b> - not enough <b>*capital*</b> to <b>*finance*</b> its early losses, which were much higher than <b>*IND</b>
3	515.407 178	8) hours a year more than the industrial average. <b>*Employees*</b> of the <b>*major*</b> oil companies form only a quarter of the workforce. The rest of the hard <b>*labour*</b> <b>*IND</b>
3	515.272 173	9) for current <b>work</b> , the current work itself, <b>*accounts*</b> , expenses and any correspondence. Clicking on items in the folders automatically loads the <b>*IND Computers:</b>

3	515.207	173	10)	throughout england and wales in a variety of <b>*occupations*</b> including environmental <b>*work*</b> , youth, arts, marketing and <b>*finance*</b> . Contact: Come along to <b>* INDEnterprise</b>
3	515.207	173	11)	high and must be reduced, and mass reduction of <b>*employment*</b> in the hitherto most stable sectors of tertiary <b>*occupations*</b> public <b>*employment*</b> , banking and <b>* INS</b> Age of extremes:
3	515.207	173	12)	in the hitherto most stable sectors of tertiary <b>*occupations*</b> public employment, banking and <b>*finance*</b> , office- <b>*work*</b> became common. <b>* INS</b> Age of extremes:

**Figure 5: Comparable contexts for "lavoro" in Italian/British newspapers**

Figure 4 gives the list of significant collocates for "lavoro" in the Italian newspapers, and Figure 5 shows the results obtained when we searched for contexts in the English newspaper archives containing a significant number of the L2 vocabulary items which were derived on the basis of this list. In Figure 5, direct translations of the term searched are shown in bold, members of the L2 vocabulary are indicated between asterisks.

### Evaluating Our Results

The comparable corpus query system is still under development, and it is clearly more difficult to evaluate its performance than it is for the parallel system. While it is easy to check whether contexts have been constructed for each direct translation derived from our bilingual dictionary for the L1 term searched, an objective evaluation of contexts that contain no direct translation but just a relevant number of items from the L2 vocabulary is not so easy. It would be necessary to go manually through the entire L2 set of texts looking for other contexts that reflect the same concept but which were not retrieved by the system to assess with some degree of accuracy where it has failed. This would be a difficult and time-consuming task, and one which we have not attempted so far. However, one test that we do make to evaluate system performance is to construct the L2 vocabulary excluding direct translations of the L1 term of interest (T). We then retrieve our comparable contexts and look to see if any of these do contain direct translations of T, despite the fact that these were not searched specifically. An example is shown in Figure 6.

This figure shows the results of a query on our comparable corpus for the Italian lemma *libertà*, using the L1 vocabulary shown in Figure 2. We show here the first 12 contexts, i.e. those calculated by the system as being most representative of the use of this term in this particular corpus. We excluded the translations of the term given in the bilingual dictionary (*liberty, freedom*) from the construction of these contexts. The fact that a direct translation of *libertà* appears in a number of the results (nos. 2, 3, 4, 5, 11 and 12) is encouraging.

### Discussion and Prospects

This approach to the problem of identifying cross-language lexical equivalences over homogeneous sets of texts for different languages has several merits: it allows us to disambiguate, to a considerable extent, both the L1 term being analysed and the target language terms provided by the dictionary; it permits us to retrieve lexically equivalent cross-language expressions even when the L2 context does not contain a dictionary-derived translation of the L1 term; it provides a ranking of four results.

DBT - (Comparable Texts) Search <i>libertà</i>				
4	24.842	20	1)	to make changes in funding <b>*policy*</b> in relation to <b>*cities*</b> , while still <b>*respecting*</b> the <b>*principle*</b> of =FE"FXAC93207ENC.0014.01.00".14
3	26.512	36	2)	Respect for human rights and <b>*fundamental*</b> freedoms is <b>*guaranteed*</b> in the Member States by <b>*effective*</b> systems of =FE "FXAC93145ENC.0030.01.00".32
3	24.462	38	3)	as long as that legislation <b>*guarantees*</b> that the relevant <b>*fundamental*</b> <b>*principles*</b> and freedoms enshrined in the =FE"FXAC93288ENC.0024.01.00".27
3	23.830	38	4)	as guardian of the Treaties. <b>*Respect*</b> for human rights and <b>*fundamental*</b> freedoms is <b>*guaranteed*</b> in the Member =FE"FXAC93145ENC.0030.01.00".32
3	22.806	40	5)	in the Treaty, such as the <b>*principle*</b> of non-discrimination, <b>*respect*</b> for the <b>*fundamental*</b> freedoms enshrined =FE "FXAC93040ENC.0013.01.00".50
3	22.806	40	6)	that all Community law ought to <b>*respect*</b> this <b>*principle*</b> because it is a <b>*fundamental*</b> right. Despite the =FE "FXAC93162ENC.0013.01.00".40



3	22.806	40	7)	countries must be based on such*principles* as *respect* for international law, human rights and*fundamental*	=FE "FXAC93327ENC.0036.02.00".26
3	22.741	41	8)	in the drawing up of all *policies* isa *fundamental* *principle* of Community legislation. In view of this: 1.=FE "FXAC93065ENC.0034.01.00".16	
3	22.741	41	9)	in Portugal. This discriminatory*policy* is in clear breach of the most *basic* *principles* and provisionsof national	=FE "FXAC93095ENC.0014.03.00".19
3	22.512	40	10)	nationality which conflict with the*fundamental* *principles* of the EEC *Treaty*. It should be made clearthat	=FE "FXAC93047ENC.0005.01.00".27
3	22.512	40	11)	under Article 115 of the EEC *Treaty*derogate from the EEC Treaty's *basic* *principle* of freedom of movementof	=FE "FXAC93264ENC.0031.01.00".32
3	22.512	40	12)	guarantees that the relevant *fundamental* *principles* andfreedoms enshrined in the EEC *Treaty* will be observed.	=FE"FXAC93288ENC.0024.01.00".27

**Figure 6: Comparable contexts for *libertà* inparliamentary texts**

**L1 or Monolingual Sense Disambiguation** Although the problem of polysemy is greatly reduced in a domain specific corpus, it is still present - to a varying degree depending on the type of texts being treated. The construction of the L1 vocabulary which characterises our term T will permit us to obtain a clustering of the most relevant terms connected to T. If the corpus contains a predominant sense for the term then the vocabulary should represent this sense - secondary senses that appear rarely will not cause a representative vocabulary of collocates to be constructed. If, in the corpus, there is more than one relevant sense for T then we would expect two or more distinct clusterings of significant collocates. Take, to use a classical example, the unlikely event that our collection of texts has a significant number of occurrences of both the river and the financial sense of "bank". We would expect to be able to obtain two distinct clusterings of significant collocates with - in this extreme case - little or no overlap. This type of sense disambiguation for the L1 term under exam did not appear very relevant when we started this work as our initial interest was in domain-specific sets of texts in different languages. However, we have now begun to extend the area of our interest to more general (although still comparable) collections, such as newspaper archives in more than one language for the same period (see above) or searching over Web sites containing documents in different languages (see final section). Being able to perform some kind of sense disambiguation on the L1 term is thus becoming far more important and we intend to pursue this line of investigation in our next work, both for monolingual and bilingual text sense disambiguation.

**L2 or Target Term Disambiguation.** The second kind of disambiguation operates at the target language level. As stated above, our procedure takes as input all the translation equivalents listed in the bilingual dictionary regardless of sense distinctions. Inappropriate translations are eliminated by the fact that we normally do not find them together with a significant number of items from the L2 vocabulary for the term being searched. For example, if we examine all the occurrences of *sicurezza* in our parliamentary corpus we find that the sense is that of "safety", or "security" (one sense of "security" is a synonym of "safety"). This is confirmed by the set of significant collocates for this term; the top ten are the Italian equivalents of toy, hygiene, reactor, health, nuclear, maritime, council, road, provisions, Euratom. The bilingual dictionary gives us four separate senses for *sicurezza*: translated by safety, security, certainty, confidence. On the English side of the corpus, we find 17 occurrences of "confidence" and just one of "certainty". However, the context for "certainty" does not appear in the list of comparable contexts for *sicurezza* as it contains no other L2 vocabulary items; the contexts for "confidence" are ranked very low as they never contain more than two L2 significant collocates for *sicurezza*. Thus, our approach helps us to identify the correct sense of the target terms offered by the bilingual dictionary and to provide a ranking of the best L2 matches for the L1 term searched. That this is not always successful, however, is shown by context No. 3 in Figure 4, in which the political sense of Labour appears although it is the work sense represented by the Italian "lavoro" that we are looking for. The reason is the very high (but unsurprising) MI value assigned to the collocates "trade" and "union"; perhaps this result is not too discouraging, after all Labour was originally the party of the workers!

The success of this approach depends on the degree to which the L1 and L2 sets of significant collocates are truly representative of the term queried. We are thus now studying ways to optimise the construction of this vocabulary. As has been stated, so far we have used the Mutual Information formula to compute our significant set of collocates for terms searched. This formula has been criticised as it tends to assign oversignificant values to infrequent words. We are currently implementing a different measure based on likelihood ratios (see

Dunning 1993). But it is too soon yet to judge whether this will give us an improvement in performance.

Our next tests should be made on (i) a more technical comparable corpus which should provide us with a real test-bed for multilingual terminology extraction; (ii) on a set of Italian and US newspaper items for the same period but of far larger dimensions to give us a chance to study the results of L1 term sense disambiguation.

## **Applications to CLIR**

In CLIR the aim is to find methods which successfully match queries formulated in one language against documents stored in other languages. Various approaches have been/are being experimented. The best-known and tested involve multilingual thesauri. Other approaches attempt to use different kinds of translation techniques in order to extend the potential range of the search: full-text translation of documents is not currently viewed as a realistic answer in consideration of the actual costs and limitations of MT systems; experiments on the automatic translation of queries employing bilingual/multilingual dictionaries have not given satisfactory results - queries are generally too short to permit an exhaustive non-ambiguous translation; likewise, concept-based approaches which attempt to achieve a matching between documents and queries at a more abstract level have not yet provided promising results on a large scale. At the moment, it looks like the most promising results will come from an integration between multilingual lexicon and corpus-based methods.

When we began work on our two bilingual corpus processing systems, our main interests were linguistic/lexicographic: applications such as bilingual lexicography, translation and language learning activities. However, we now recognise that both systems, employed together with our bilingual electronic lexicons, can also be applied to CLIR activities.

### ***Parallel Text System***

The very "tight" alignment achieved using dictionary derived translation links would greatly facilitate the statistical alignment of unmatched terms. So far, this system has been tested in human-oriented applications, we now want to experiment it over a large parallel collection in order to automatically extract "new" translation equivalent data and thus augment the existing bilingual lexicon.

### ***Comparable Text System***

We intend to test our comparable-corpus strategy on two applications:

- multilingual digital library
- multilingual web search engine

In both cases, we will integrate a dictionary/thesaurus-based search with corpus-based strategies. Disadvantages of lexicon-based systems are that thesauri confine users to a controlled-vocabulary while general-purpose dictionaries tend to be lacking in necessary technical vocabulary; the problem with most corpus-based CLIR systems is that the acquisition of a suitable set of relevant documents on which to train the retrieval system is very resource consuming. We hope to overcome these two problems: the comparable-corpus strategy can be used to extend the limits of a simple query term translation, and also to reduce the risk of ambiguity in the query term, and to provide a ranking of the results; at the same time, our corpus will consist of the documents in the collection being queried and does not have to be created ad hoc.

Queries will be translated by the multilingual lexicon but will also be expanded by applying the comparable-corpus based strategy in order to associate with each query term, not only its direct translations but also a vocabulary which defines its probable immediate context, in L1 and L2. In this way, we search for both pre-identified translation equivalents and also cross-language lexical equivalences. When the dictionary or lexicon offers no translation equivalent, the search for cross-language equivalent contexts is still possible. Documents retrieved are ranked with respect to (i) translation equivalents of query terms, (ii) statistical value assigned to associated significant collocates.

## References

- Brown, P.F., Della Pietra, S.A., Della Pietra, V.J., Mercer, R.L. (1993). The Mathematics of Statistical Machine Translation: Parameter Estimation. *Computational Linguistics*. 19(2):263-312.
- Church, K.W., Hanks, P. (1990). Word Association Norms, Mutual Information and Lexicography. *Computational Linguistics*. 16(1): 22-29.
- Dunning T. (1993). Accurate Methods for the Statistics of Surprise and Coincidence, *Computational Linguistics*, 19(1).
- Gale, W.A., Church, K.W. (1993). A program for aligning sentences in bilingual corpora. *Computational Linguistics*, 19(1): 75-102.
- Hartmann, R.R.K. (1997). From Contrastive Textology to Parallel Text Corpora: Theory and Applications. In Raymond Hickey and Stanislav Puppel (eds.), *Language History and Language Modelling*. Festschrift in Honor of Jacek Fisiak's 60th Birthday, Berlin: W. de Gruyter, in press.
- Marinai, E., Peters, C., Picchi, E. (1991). Bilingual Reference Corpora: A System for Parallel Text Retrieval. In *Using Corpora*, Proc. of 7th Annual Conference of the UW Centre for the New OED and Text Research. Oxford: OUP, 63-70.
- Marinai, E., Peters, C., Picchi, E. (1992). A Project for Bilingual Reference Corpora. *Acta Linguistica Hungarica*, 41 (1—2). Akadémiai Kiadó, Budapest, 1-15.
- McEnery, T., Oakes, M. (1996). Sentence and word alignment in the CRATER project. In J. Thomas and M. Short (eds.), *Using Corpora for Language Research*, Longman, London and New York, 211-231
- Picchi, E., Peters C. (1996). Cross Language Information Retrieval: A System for Comparable Corpus Querying. In G. Grefenstette (ed.) Working Notes of the Workshop on Cross-Linguistic Information Retrieval, ACMSIGIR'96, 24-33.
- Smeaton, A.F. (1992). Progress in the Application of Natural Language Processing to Information Retrieval Tasks, *The Computer Journal*, 35 (3), 268-278.
- Warwick-Armstrong, S., Russell, G. (1990). Bilingual Concordancing and Bilingual Lexicography. In *EURALEX 4th International Conference*, Malaga, Spain.