



Integrating Corpus-based Resources and Natural Language Processing Tools into CALL

PASCUAL CANTOS-GOMEZ*
Universidad de Murcia

ABSTRACT

This paper aims at presenting a survey of computational linguistic tools presently available but whose potential has been neither fully considered nor exploited to its full in modern CALL.

It starts with a discussion on the rationale of DDL to language learning, presenting typical DDL-activities, DDL-software and potential extensions of non-typical DDL-software (electronic dictionaries and electronic dictionary facilities) to DDL.

An extended section is devoted to describe NLP-technology and how it can be integrated into CALL, within already existing software or as stand alone resources. A range of NLP-tools is presented (MT programs, taggers, lemmatizers, parsers and speech technologies) with special emphasis on tagged concordancing.

The paper finishes with a number of reflections and ideas on how language technologies can be used efficiently within the language learning context and how extensive exploration and integration of these technologies might change and extend both modern CALL and the present language learning paradigm.

KEYWORDS: Concordancing, corpus linguistics, data-driven learning, electronic dictionaries, modern language technologies, linguistic corpora, machine translation, morphological generation, NLP-tools, parsing, POS-tagging, speech technology

* *Address for correspondence:* Pascual Cantos Cóñez, Departamento de Filología Inglesa, Universidad de Murcia. C/. Santo Cristo 1, 30071 Murcia, Spain, e-mail: pcantos@um.es.

I. INTRODUCTION

The use of concordancing in literature and linguistic analysis is nothing new. It started well before computers existed. Tribble and Jones (1990) trace the history of concordancing from the 13th century, when Hugo de San Charo enlisted 500 monks in producing a complete concordance of the Latin Bible. That is, a kind of reference work designed to assist in the exegesis of the Bible, consisting of all occurrences of terms, names, etc., that were felt to be significant, and presented these terms in a way that would help the researcher.

Of course, current applications of concordances in language and literature are not so labour intensive. The use of concordancing as a tool for language learning/teaching is relatively recent, starting in the 1980's, when computational power began to get scaled into small, affordable personal computers.

Succinctly, a concordance is a data arrangement technique that transforms texts into lists, printing lines of text where the word or expression interested in investigating is displayed in the centre of line, know as *KWIC* (Key Word In Context), within an arbitrarily selected context of characters or words to its right and left. This technology permits that, for example, a language teacher or learner interested in knowing the use of the preposition *of* to transform a text such as:

```

What be more important or intriguing than our own origins? Like all animals
we come from one cell that develops into an embryo which forms the adult.
This embryonic development presents a fundamental problem of biological
organization. From the single cell, the fertilized egg, come large numbers
of cells -many millions in humans- that consistently give rise to the
structures of the body. How do these multitudes of cells become organized
into the structures of, for example, our body -nose, eyes, limbs, and brain?
What controls their individual behaviour so that a global pattern emerges?
And how are the organizing principles, as it were, embedded or encoded
within the egg? It is remarkable that a cell as overtly dull and
structureless as the fertilized egg can give rise to such varied and complex
forms. The answer lies in cell behaviour and how this behaviour is
controlled by genes.
  
```

Into the following format:

```

... presents a fundamental problem [[of]] biological organization.
... egg, come large numbers [[of]] cells -many millions in ...
... rise to the structures [[of]] the body. How do ...
... How do these multitudes [[of]] cells become organized into...
... organized into the structures [[of]], for example, our body ...
  
```

Clearly, what this technique does is to make the invisible visible for teachers and learners. Patterns and structures that would else hardly be immediately recognisable, spring to the eye.

II. A WORD ON DATA-DRIVEN LEARNING

There is ample discussion in the literature on the merits of linguistic corpora in second language teaching and learning (Aston 1995, 1996, 1997; Ball 1995; Barlow 1992, 1995; Burnard and McEnery 1999; Celce-Murcia 1990; Collins 1999; Flowerdew 1999; Gavioli 1997; Higgins 1991a, 1991b; Johns 1986, 1988, 1991, 1993; Johns and King 1991; Leech and Candlin 1986, etc.), mainly as a result of the pioneering work of Tim Johns at Birmingham University (1986, 1988).

Johns (1988) states that the use of concordancing in language learning: (a) interjects authenticity (of text, purpose, and activity) into the learning process; (b) learners assume control of that process; and (c) the predominant metaphor for learning becomes the research metaphor, as embodied in the concept of data-driven learning (DDL), which builds learners' competence by giving them access to the facts of linguistic performance.

Higgins (1988) proposes *concordances* as the central idea to shift the pedagogical teaching/learning paradigm from computer as magister to computer as pedagogue. That is, from mere process-control model of language instruction to an information-resource model, where learners explore, hypothesize and learn the language for themselves and the role of instruction is to provide tools (concordance programs) and resources (texts or corpora) for doing so. Similarly, Cobb (1997) considers that DDL has a specific learning effect that can be attributed to the use of concordance software by language learners. He concludes that computer concordances might simulate and potentially rationalize off-line vocabulary acquisition by presenting new words in several contexts.

Stevens (1995) accounts that many teachers feel that concordancers are the type of software that most closely approaches fulfilling the potential of computers in language learning. In a sense, they are working approximations of expert systems. They bring cognitive and analytic skills in students to bear on the manipulation of comprehensive databases for the purpose of solving real-language problems.

The effectiveness of concordances becomes also apparent not just in teaching/learning, but also in linguistic research. By means of this technique, Kettemann (1995) compares the treatment of the English conditional clause in a standard grammar used in Austria (Kacowsky 1987) with the evidence of authentic usage (corpora), and his comparison showed that an important type of English conditional with present tense in both main clause and conditional clause, accounting for one third of all instances in the data, is ignored in Kacowsky's grammar.

Tribble (1997) stresses the potential and usefulness of DDL to language learning/teaching even with few corpus resources or small, specific corpora. Corpora and corpus-based exercises are useful because they favour learning by discovery —grammar, vocabulary, etc.— (Tribble and Johns 1990:12).

A further related issue with DDL is authenticity. Widdowson (1983:30) considers that

An authentic stimulus in the form of attested instances of language does not guarantee an authentic response in the form of appropriate language activity [...] we should retain the term 'authenticity' to refer to activity (i.e. process) and use the term 'genuine' to refer to attested instances of language.

In this sense a corpus may contain millions of "attested instances of language". but there is nothing to guarantee that you can use data from that corpus as a stimulus for "appropriate language activity" (Tribble 1997). That is, it is likely that foreign language students are not necessarily motivated by a language learning activity if the instances of language use that they are studying are extracted from contexts that have little or no connection with their interests and concerns. Genuine examples of language in use will not necessarily lead to authentic language use or effective language learning activities.

So the question is: which is the best corpus for language learners? Flowerdew (1993:309) thinks that

Many native speakers make use of others' writing or speech to model their own work in their native language where the genre is unfamiliar. It is true that this skill was brought out of the closet, and exploited as an aid for learning.

Similarly, Bazerman (1994:131) considers that the most useful corpus for learners of English is the one which offers a collection of expert performances in genres which have relevance to the needs and interests of the learners. These texts might exemplify the results and models of the desired forms of language behaviour that language learners want to achieve and might, therefore, be motivating starting points for language learning and language using activities.

Clearly, this, somehow, relegates standard, balanced and representative corpora, such as the Brown corpus of American English (Kučera and Francis 1967), the Lancaster-Oslo-Bergen (LOB) corpus of British texts (Johansson 1980), or other major corpora such as the British National Corpus (BNC) (Burnard 1995), a 100 million word representative corpus of contemporary British written and spoken texts, or the Bank of English at Birmingham University (Sinclair 1991), for language learning purposes. Tribble (1997) points towards non-standard corpora for DDL and draws his attention to multimedia encyclopaedia, such as Microsoft Encarta®, among others. The latest version of Microsoft Encarta® contains more than 30,000 articles, between 200 and 5000 words, which count for a total of roughly 30 million words, covering different domains and topics, such as art, geography, history, language, life science, literature, philosophy, physical science, religion, social science, sports, etc. The data provided by this multimedia encyclopaedia virtually contain enough texts which most students in most language classes will find interesting and informative. In addition, with this comprehensive range of topics and texts, it is not difficult to select *ad hoc* texts, focussing on students particular needs and motivations.

Of course, our aim here is not to advertise any particular multimedia encyclopaedia but much more to encourage language students and teachers to use the vast range of language texts, corpora or data, in general, which is available in electronic form, in CD-ROMs and/or Internet, rather than urging them to construct their own comprehensive and representative corpora.

II. A WORK ON EXISTING DDL-SOFTWARE

In this section, we shall review the main software applications used among DDL practitioners: (1) commercially available concordance programs and (2) Tim Johns' *Context*. In addition, we shall also present a Spanish vocabulary learning multimedia application, *Practicatu Vocabulario* (Sánchez and Cantos 2000), which is based on the electronic dictionary metaphor and DDL-like learning/acquisition strategies.

II.1. Concordancers

Concordancers are text processing tools for looking at how words behave in texts. These tools allow you to find out how words are used in texts. Among the facilities, all concordancers allow you at least¹:

- To list all the words or word-clusters in a text, set out in alphabetical or frequency order.
- To see any word or phrase in context (concordances), so that you can see what sort of company it keeps.

This text processing tool is generally used for lexicographic work, for preparing dictionaries, and by researchers investigating language patterns.

Tim Johns has compiled numerous exercise examples on his web page² using standard concordancing tools. The classroom materials that follow are extracts from his website and are a collection of some participants' work of the *Ustí nad Labem DDL Workshop* (21st - 25th March 2000):

About and On

How do we use the prepositions ABOUT and ON? Which one is used more often? Which one tends to be used in academic texts? In which cases is there the occurrence of only one of them?

Book

I have sent a British Medical Association book about a potential risk to human health to the market. We published a coffee-table book about ant behaviour called The Ant, which costs sterling 16.95. Publicised as a book about the terrible fate awaiting humanity in the future. This is yet another 'gee-whizz' book about forensic science, this time based on a true story. I was intrigued by the title: at last, I thought, a book about the personal relationships that science

has seen evergreen. I remember buying my first book on planets (by Patrick Moore) back in the 1960s. How things have changed. A good book on the planets has always needed to be updated. Afghanistan, so how did he come to write a book on Murdock? Was his choice dictated by the 24 hours? I'm keeping that for myself. It's a book on kilims - geometrically-patterned rugs from Persia (HarperCollins, 1990), and is writing a book on the future of US national security policy.

(by Květa Rychtářová)

'Great', 'Big', 'Large', 'Huge' and their Collocations

Task 1 - 'Great' and its collocations

Try to spot what the typical cases of 'great' and its collocations are on the basis of the following examples.

- Whether recent discoveries in the Great Pyramid of Cheops have anything to do with the discovery of the tomb of the pharaohs is a matter of debate. The great work of restoring Al-Andalus and its Great Mosque to their former faith is a task that has occupied many generations. The great king himself was comfortable. And the great London clubs, with their roaring coal fires and smoking pipes, were a part of the life of the city.
- The great narratives of Les Misérables and Great Expectations. What could be more comical than the story of a man who had been a great success in the Great Journey series, the theme of which was the search for a better life.

Task 3 - 'Great' and missing nouns

Try to predict the words which are missing.

- Falstaff who has been lured into Windsor Great _____ with antlers on his head.
- Benazir is one of Europe's great _____, like Rodbet an all-round.
- _____ had become known as Great _____ The Times, the Daily Mail and the Daily Express.
- In 1679 Frederick II the Great, _____ of Prussia.
- The world was dominated by the great European _____ and, since the 1850s, by the United States.

(Park, players, race, King, Empires)

(by Sarka Canova and Jarka Ivanova)

Adjectives ending in -ic, -ical

I. Look carefully at the following citations. What difference in the meaning of the adjectives *classic/classical* can you spot?

And when Sinatra was making his classic albums for Capitol in the 1950s *Songs For Swinging ...* Gramophone, the classical music magazine, has not written about her single or her album.

II. How fill in the blanks with appropriate adjectives.

1. I remember listening to all the _____ Motown music and the Philadelphia soul stuff.
2. Once it became clear that she could not continue through the next two acts, Deane decided to replace both dancers. As in _____ ballet one partner may not be physically suited to perform with a stand-in.
3. ... we always dine at Cafe Des Arts. A _____ bistro, run by a consortium of charming ladies, stylish, innovative food.
4. The area where the dirt collects is transparent. All our detritus is paraded on the outside, turning the _____ design inside out. Why do we need to see it?
5. ... architect Richard Norman Shaw, capable of turning out gothic, Queen Anne and strict _____ designs, made a valuable contribution to the Arts and Crafts movement.

(by Zuzana Šaffková and Vladislav Smolka)

II.2. Ready-made DDL-software: *Context*

Probably one of the best known and most used DDL-software is *Context*³. This program encourages language learner to investigate how words are used in context in English, and is designed to supplement classes. It is based on short contexts (extracted from the database by means of the computer program *MicroConcord*) illustrating the use of important key items from a database of over 3 million words of text in English.

The program starts offering the user a list of headings: *Top Menu* (parts-of-speech and topics; see *Figure 1*). In addition, it is also possible to view a more detailed index of all the keywords available to the program, together with the names of the files in which each key word is stored and to select a file of contexts (by keywords defined by parts-of-speech, keywords defined by topic or morphemes—prefixes or suffixes; *Figure 2*). Once the user has selected the file of context, the program displays the list of key items in the bottom of the screen (among other facilities) in order to investigate the set of contexts for any particular key item (*Figure 3*). The *Quiz Screen* challenges students to guess what the missing keyword is (*Figures 4 and 5*). After students have finished the Quiz, they can see an analysis of their performance.

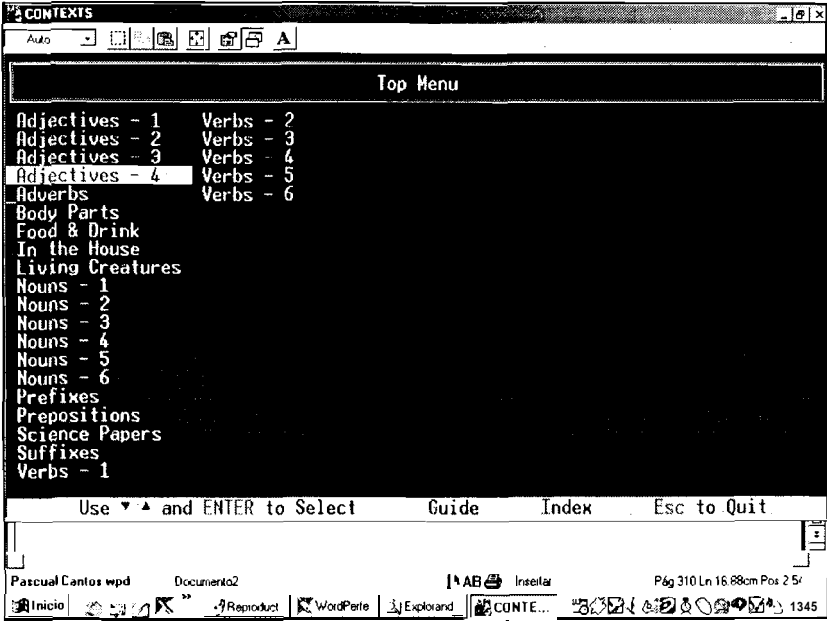


Figure 1. Top Menu

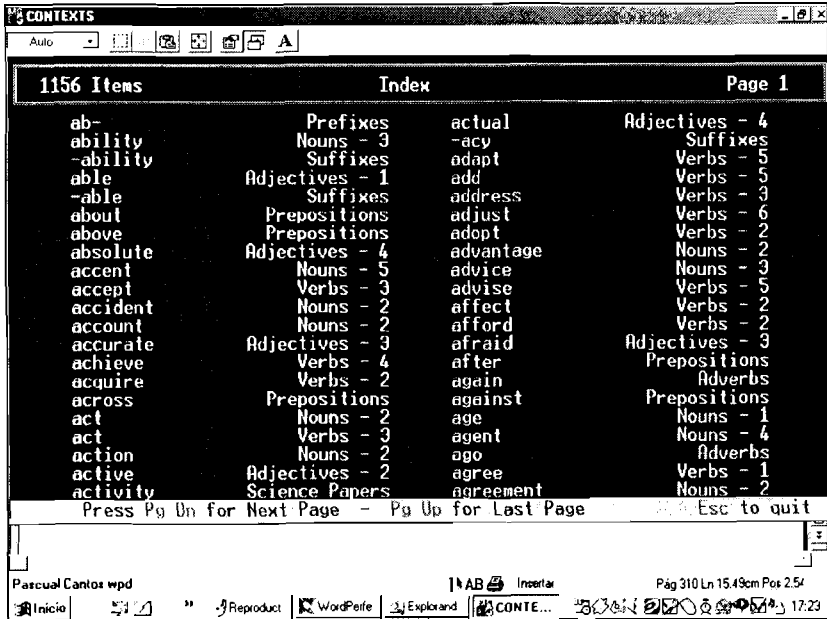


Figure 2. Indexed-data Window

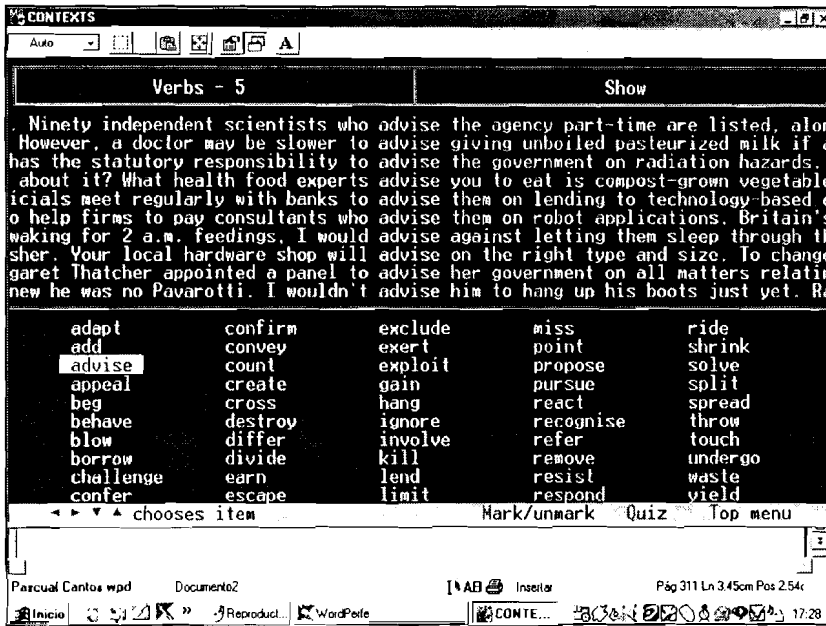


Figure 3. Concordance-data Window

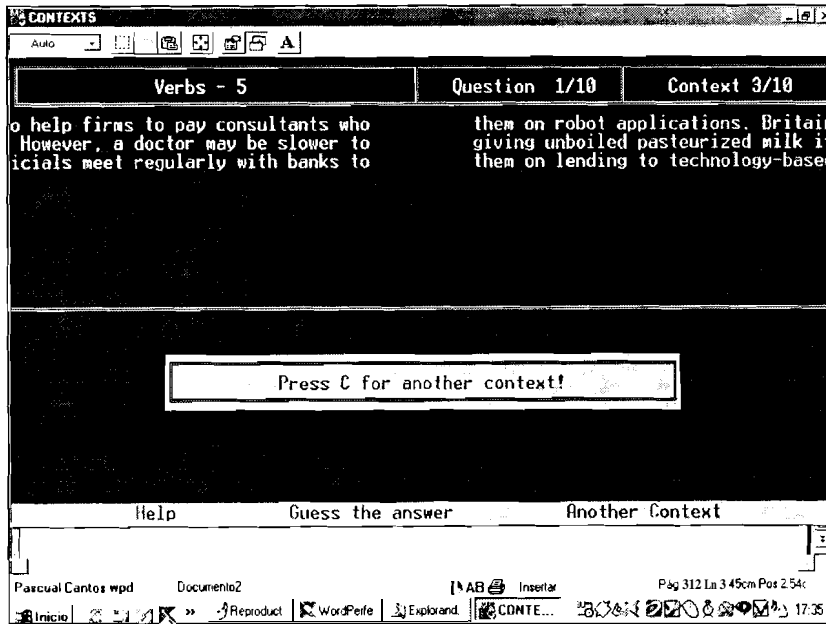


Figure J. Quiz Window (1)

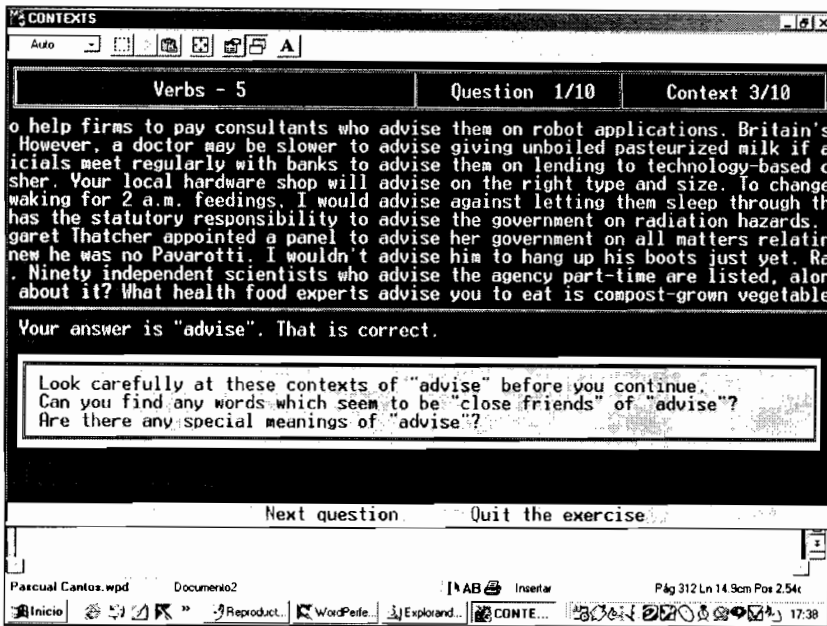


Figure 5. Quiz Window (2)

II.3. Using Electronic Dictionary Facilities for DDL

Succinctly, electronic dictionaries are commercial derived products of standard paper dictionaries. The main differences between paper and electronic dictionaries is that in the latter you can find words immediately, even if you are unsure of the exact spelling. They can also be run as a stand-alone program and can be used in conjunction with most word-processing software; you can browse through entries, view adjacent entries, or travel swiftly between entries. In addition, some electronic dictionaries keep track of your searches so that you can very easily return to words you have previously looked up. You can also print extracts or definitions and copy them to the clipboard. Electronic dictionaries can be monolingual, bilingual or multilingual (Figures 6 and 7).

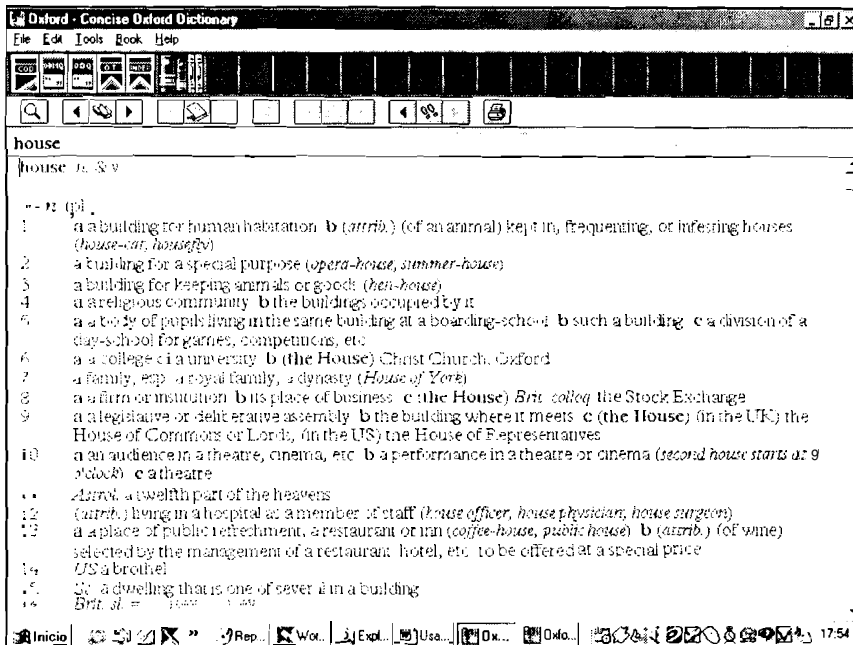


Figure 6. Example of a monolingual electronic dictionary

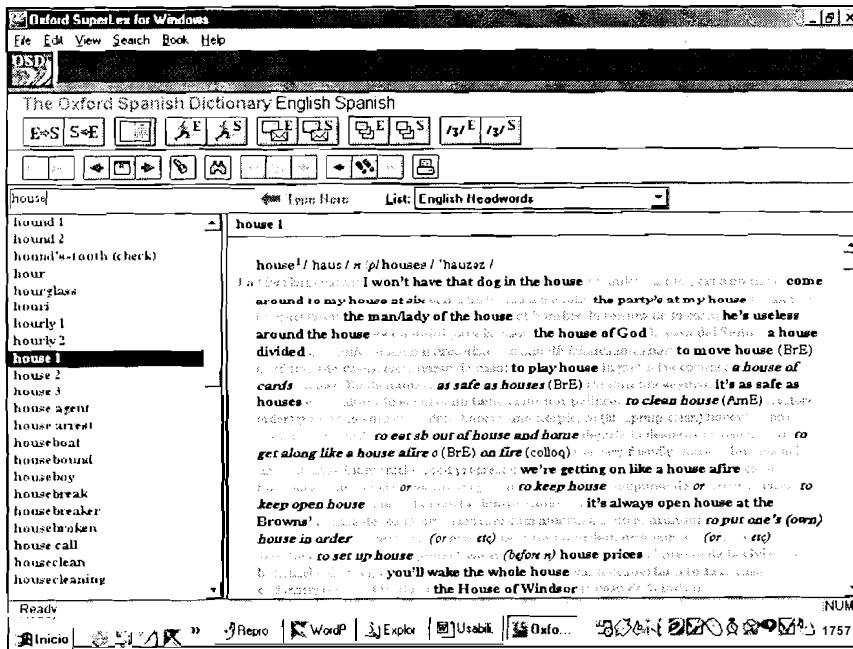


Figure 7. Example of a bilingual dictionary

Based on the electronic dictionary metaphor, Sánchez and Cantos⁴ designed a DDL-like software: *Practica tu Vocabulario (PTV)*. *PTV* is a Spanish lexicon learning software, containing the 4500 most frequent types occurring in the *CUMBRE Corpus* —a linguistic corpus of contemporary Spanish (Sánchez et al. 1995). All 4500 items:

Are translated into English, French, German, Portuguese and Italian. By just clicking on the desired flag, students will get the words translated in that language. However, students might change translation language any time at will (Figure 8).

Can be accessed, using standard electronic search facilities: term search, window scroll or thumb index (Figure 9).

- Are illustrated with a real example —full concordance sentence, extracted from the *CUMBRE Corpus* (Figure 10).
- Are recorded and can be heard by the students.



Figure 8. Language Selection Window

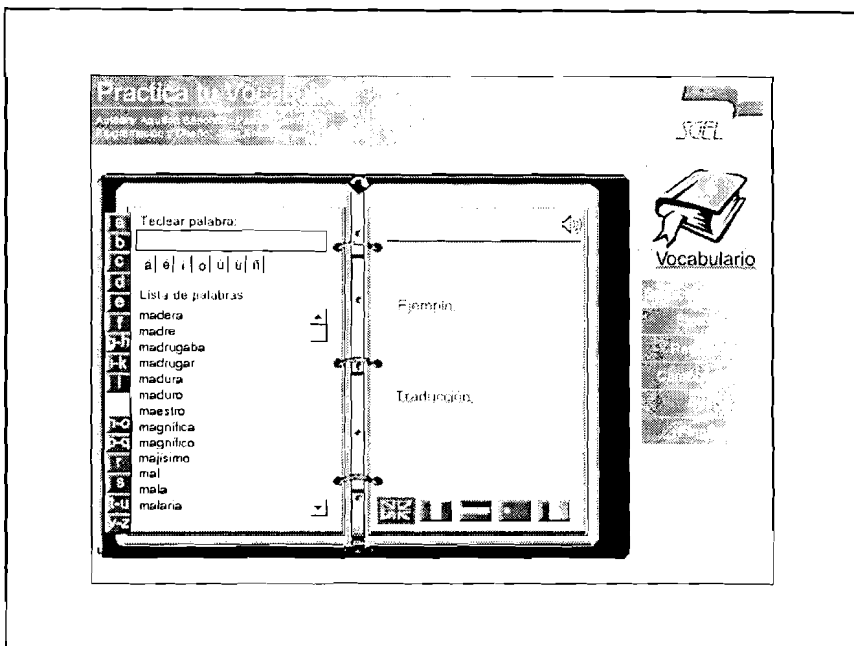


Figure 9. Search Facilities

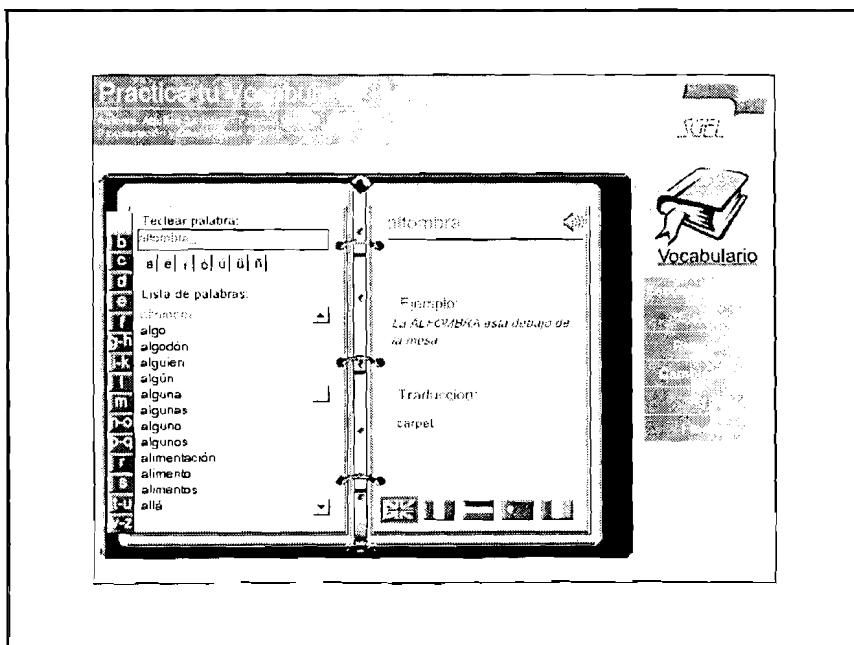


Figure 10. Visualising Concordance-sentence and Translated Term

The option *Ejercicios* offers three types of exercises:

Listen, repeat, record and check your pronunciation. On selecting this exercise, students will choose the number of words they wish to work with by clicking on the button with the number of examples which will, at random, be the basis of the exercise. Next, the random-selected words will appear, and by clicking on the loudspeaker icon, the blue-highlighted word can be heard. Finally, students click on the microphone icon and record the highlighted word. A click on the right loudspeaker reproduces the *model* recording followed by the student's recording and the student can contrast both outputs (*Figure 11*).

Listen and write. This is a word dictation practice; students will hear randomly chosen words and will have to write them correctly. The program allows three guesses before displaying the correct spelling. To facilitate the writing of Spanish diacritics, PTV provides them on a small table below the text-entry window (*Figure 12*).

Read in your language and translate into Spanish. Here the program displays randomly selected words in the target language chosen and students have to write the translation for each word into Spanish (*Figure 13*).

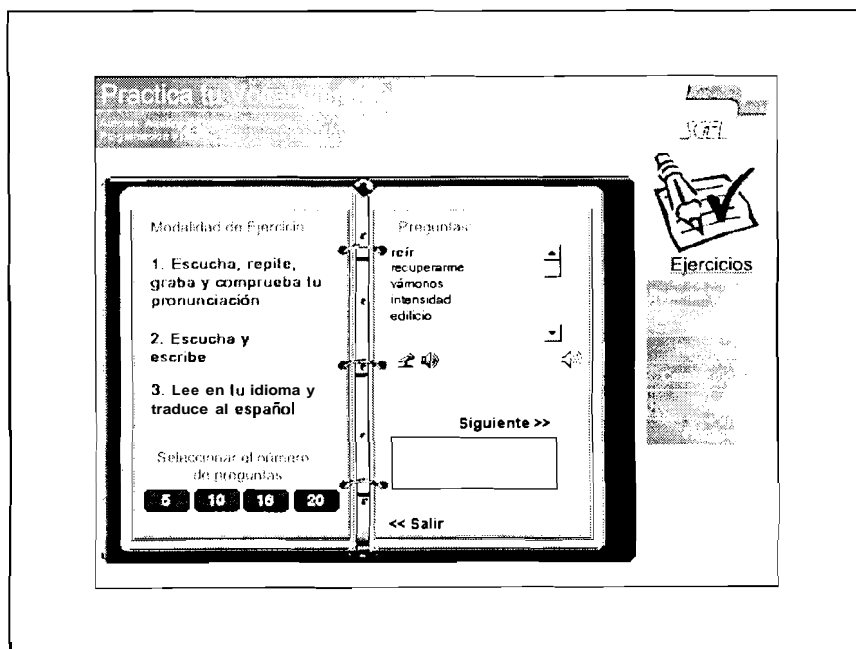


Figure 11. Exercise Type 1: *Listen, Repeat, Record and Check your Pronunciation*

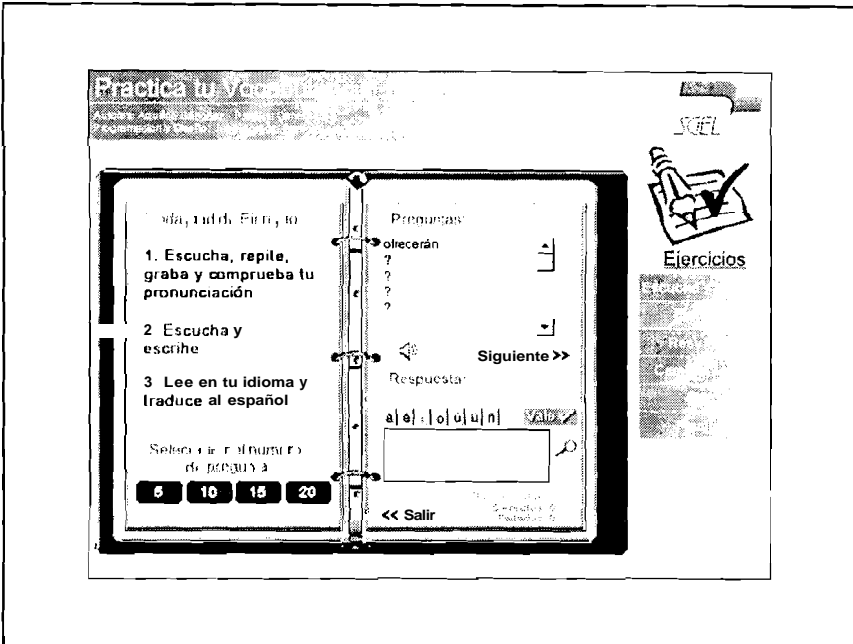


Figure 12. Exercise Type 2: Listen and Write

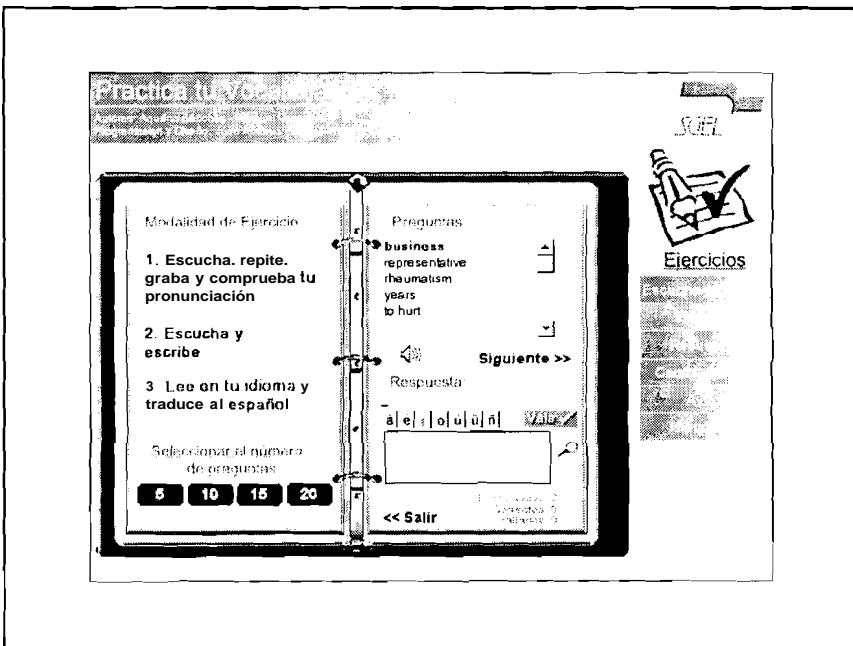


Figure 13. Exercise Type 3: Read in your Language and Translate into Spanish

On completion of each exercise or at the end of a working session, students may consult his/her lists or failures by clicking on the *Resultados* tab on the right of the agenda.

III. INTEGRATING DDL AND LANGUAGE TECHNOLOGIES

In recent years, a new term has been coined by the CALL community: *Human Language Technologies* (HLT). This term embraces a wide range of research and development areas within the area of Language Engineering or Language Technologies.

The field of human language technology covers a broad range of activities with the eventual goal of enabling people to communicate with machines using natural communication skills. Research and development activities include the coding, recognition, interpretation, translation, and generation of language. ... Advances in human language technology offer the promise of nearly universal access to on-line information and services. Since almost everyone speaks and understands a language, the development of spoken language systems will allow the average person to interact with computers without special skills or training, using common devices such as the telephone. These systems will combine spoken language understanding and generation to allow people to interact with computers using speech to obtain information on virtually any topic, to conduct business and to communicate with each other more effectively. (Cole 1996)

III.1. Some HLT tools⁵

There are many HLT tools that have become commercial systems. Among those systems, probably the two areas that have focused most commercial and scientific motivation are Machine Translation (MT) and Speech Recognition (SR). Particularly interesting here is the possible application domain of MT and SR to CALL and more generally, language teaching and learning, and other HLT tools. Interesting in this respect are part-of-speech (POS) taggers and syntactic parsers. These two Natural Language Processing (NLP) tools might help teachers and learners to preprocess texts and highlight certain grammatical phenomena or patterns without the trouble of having to manually annotate a text.

In the following sections, we shall introduce some HLT applications and try to highlight their interest for language teachers and learners, in general, and also for non-HLT initiated CALL practitioners.

III.1.1. Machine Translation⁶

From the earliest days, MT has been bedevilled by grandiose claims and exaggerated expectations. In present day, however, the term MT is generally the standard for computerised systems responsible for the production of translations from one natural language into another, with or without human assistance.

Although the ideal may be to produce high-quality translations, in practice the output of most MT systems is revised and edited. In this respect, MT output does not differ much from the output of most human translators which is normally revised by other translators before dissemination. MT output may also serve as rough or raw translations.

While many of the commercially available MT packages may be useful for extracting the gist of a text they should not be seen as a serious replacement for the human translator. Most machine translations are not that bad, they are *half-intelligible*, letting you know whether a text is worth having translated properly and there are many situations where the ability of MT systems to produce reliable, if less than perfect, translations at high speed are valuable. Even where the quality is lower, it is often easier and cheaper to revise 'draft quality' MT output than translate it entirely by hand. The translation quality of MT systems depends mainly on restrictions of the translation domain, linguistic architecture and components.

Imposing restrictions on the input such as (a) limiting the texts to particular sublanguages of document type and subject field and/or (b) controlling the language (reducing ambiguities, colloquial expressions, etc.), may improve translation quality.

Regarding MT architecture, the first MT systems are generally referred to as having a direct translation approach. The main idea behind this architecture is that source language sentences can be transformed into target language sentences by *shallow* analysing the source text, replacing source words with their target language equivalents as specified in a bilingual dictionary, and then roughly re-arranging their order to suit the rules of the target language. The second basic type is the interlingua approach. This type assumes the possibility of converting texts to and from *meaning* representations common to more than one language. Translations consist of two stages or phases: (1) from the source language to the interlingua and (2) from the interlingua to the target language. The third type of MT systems, the transfer approach, involves three stages: (1) converting source texts into intermediate representations in which ambiguities have been resolved irrespectively of any other language, (2) converting these into equivalent representations of the target language, and (3) generating the target texts (translations).

Some other MT systems rely less on the approaches mentioned above. Example-based machine translation, for instance, does not employ mapping between languages but instead matches stored translation examples against each other using a bilingual corpus of translation pairs (Nagao 1984). An even more radical approach to MT is the statistical approach (Brown et al. 1993) which requires the use of large bilingual corpora which serve as input for a statistical

translation model.

Regarding language teaching, MT systems can be easily and efficiently integrated into the learning process. Some potential applications are⁷:

- Translating full texts or paragraphs. Student can translate and then read the texts in their own language, extracting the gist without teacher intervention (Figure 14). In most MT systems, users can translate sentences automatically or interactively. Automatic translation proceeds autonomously, without the intervention of the user, whereas in interactive translation, the user can intervene in the translation process and choose the best word whenever more than one translation is possible.
- Translating sentence-by-sentence and print the source and target texts in a line-by-line format. This layout can be useful for comparing the original and translated text. This allows students to explore for equivalents between the source and target language, look for erroneous translations/false friends and assist in their own translations. (Table 1).
- Studying or writing in a foreign language.
- Looking up words (dictionary) and their inflections (Spanish grammar) (Figure 15).

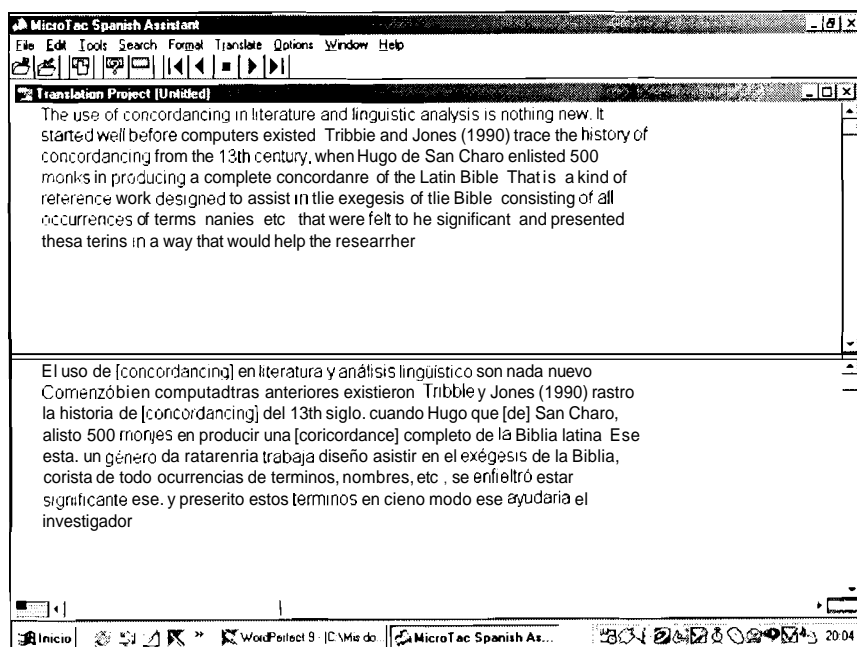


Figure 14. Automatic Text Translation (English-Spanish)

1.	Tlic use of coicordanciig in literature and linguistic analysis is nothing riew. El uso de [coicordancing] en literatura y análisis lingüístico son nada nuevo.
2.	It staited well bfore coinputers existed. Comenzó bien computadoras anteriores existieron.
3.	Tribble aiid Joies (1990) trace tlie history of concordancing from the 13th century, when Hugo dc San Charo enlisted 500 monks in producing a complete coricordance of the Latin Bible. Tribble y Joies (1990) rastro la liistoria de [concordancing] del 13th siglo, cuando Hugo que [dc] San Clario alistó 500 inonjcs en producir una [concordance] completo dc la Biblia latina.
4.	That is, a kiind of refrcrncc work designed to assist iii tlic exegesis of tlie Bible, coisistiig of all occurrciics of tcrriiis. naliies. etc., that were felt to be significant, and presented these tcrriis in a way that would help tlie rescarclir. Ésc está, un género de referencia trabaja diseñó asistir en el exégesis de la Biblia. coiista dc todo occurciias de térrriios, nombres, etc. se entielró estar significante ésc, y presciitó estos términos en cierto modo ése ayudaría el irvestigador.

Table 1. Line-by-line Printed Translation

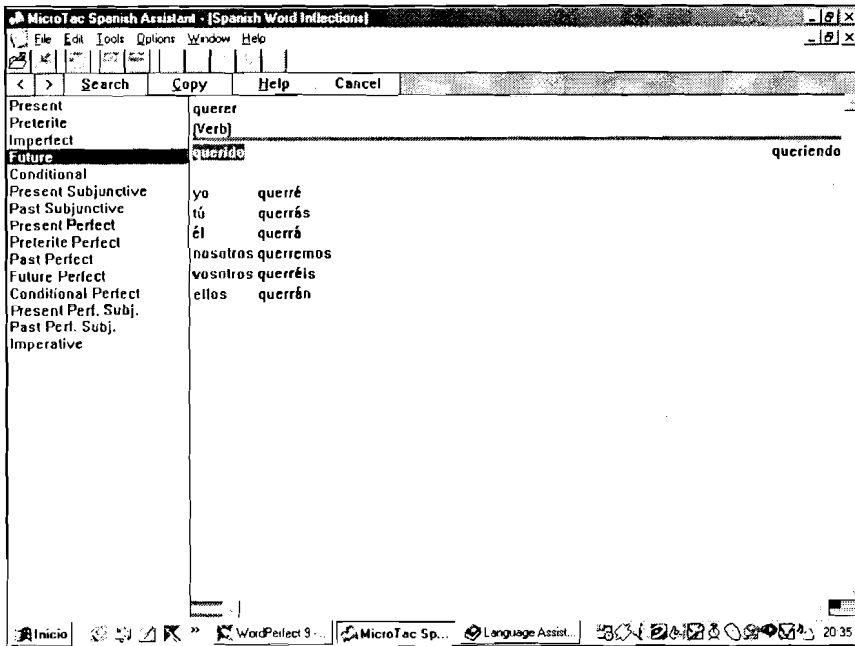


Figure 15. Intlection Look-up

III.1.2. Part-of-Speech Tagging

Word frequency lists derived by computers from corpora have clear shortcomings. These computer counts and sorting of word forms somehow bury or distort important facts about the language: variant inflected forms of nouns and verbs especially would be treated as entirely different word types. For example, *be*, *am*, *are*, *is*, *was* and *were* would be accounted for completely different linguistic items. Similarly, the frequency count of the number of occurrences of the word form *light* in a corpus would include the noun, verb, adjective and adverb.

Because manual annotation of each word token with its parts-of-speech (POS) in the corpus would be too expensive, the solution adopted has been to design computer programs, known as *POS-taggers*, to annotate automatically every word in the corpus with a *tag* to show the POS it belongs in context.

*TAGGIT*⁸ was the first computer program designed and implemented to annotate a major corpus and assigned 87 tags to the word forms in a corpus. Subsequent developments in POS-taggers found necessary to expand the tagset and to modify the rule-based approach of *TAGGIT*.

Other tagging systems, such as *CLAWS*⁹, are based on probabilistic principles and are remarkably robust. In particular, *CLAWS* uses 133 basic word and punctuation tags and gets a minor error margin of just 3-4%.

Another extension of automatic POS-tagging is the combination of rules and stochastic or probabilistic principles. This is found in *eTiKeT@*¹⁰. Actually, this HLT tool is not a tagger but a tagger-generator. It has not been designed for any specific language but, in principle, for any language. It starts from scratch: with an empty lexicon (data base), without any linguistic information (rules) nor probabilistic data and uses just 14 tags (*Table 2*). The user's task is to train or "teach" it for the language desired. All sessions are stored and the manual tagging is compared with the system's performance (*Tables 3a, 3b and 4*). Once a satisfactory success rate has been achieved, the system can be left to perform automatically without human intervention.

To speed up the initial human tagging phase, the user can alternatively feed the system's lexicon with stoplist items. That is, high-frequent non-ambiguous types, mostly close-class items, such as pronouns, prepositions, conjunctions, articles, auxiliary and modal verbs, etc.

The program tags on a sentence-by-sentence basis and outputs the results either in a database mode (*Figure 16*) or as running ASCII text with the tags attached to the tokens in the text (*Figure 17*). Additionally, the user can also consult the patterns and statistics the system has inferred so far (*Figure 18*).

ID	POS	Abbreviation
0	Noun	N
1	Verb (lexical)	V
2	Verb (aux)	Aux
3	Verb (modal)	Mod
4	Adjective	Adj
5	Adverb	Adv
6	Preposition	Pre
7	Particle	Par
8	Conjunction	Con
9	Interjection	Int
10	Determiner	Det
11	Pronoun	Pro
12	Punctuation	Pun
13	Other	Oth

Table 2. Tag-set of *eTiKeT@a*.

SessionCode	FileCode	CorrectGuess	WrongGuess	LeftContext	RightContext
2	2	2	2	3	3
3	3	3	1	3	3
4	4	5	0	3	3
5	5	6	1	3	3
6	6	6	1	3	3
7	7	2187	1407	3	3
8	10	63	50	3	3

Table 3a. Information on Session Performance and POS-Disambiguation Context selected

TagLastWord	Date	FinishedSession	JustWords	Language
0	2001-04-19 14:51:09.36	YES	NO	0
0	2001-04-19 14:52:35.21	YES	NO	0
0	2001-04-19 14:54:48.02	YES	NO	0
0	2001-04-19 14:57:17.85	YES	NO	0
0	2001-04-19 14:58:53.97	YES	NO	0
2	2001-05-10 15:25:29.61	YES	NO	0
6	2002-04-25 10:20:48.49	YES	NO	0

Table 3b. Information on Tags, Date, Session, Text or Single Word Tagging and Language¹

Word	POS	Frequency	SessionCode	CorrectGuess	Language
a	10	62	0	YES	0
about	6	17	7	YES	0
above	6	1	0	YES	0
accept	1	1	7	NO	0
acceptable	4	1	7	NO	0
achieve	1	1	7	NO	0
achieved	1	1	7	NO	0
achieves	1	1	7	NO	0
achieving	1	1	7	NO	0
across	6	1	7	YES	0
actual	4	3	7	YES	0
additional	4	2	7	NO	0
address	0	1	7	NO	0
address	1	2	7	NO	0
addressed	1	1	7	NO	0
administration	0	2	7	NO	0
admiral	0	1	7	NO	0
advance	1	2	7	NO	0
affect	1	1	7	NO	0
after	5	1	7	YES	0

Table 4. Data Base Extract (Types. Tag. Frequency. Session. First-Time Guessing of the Type and Language)

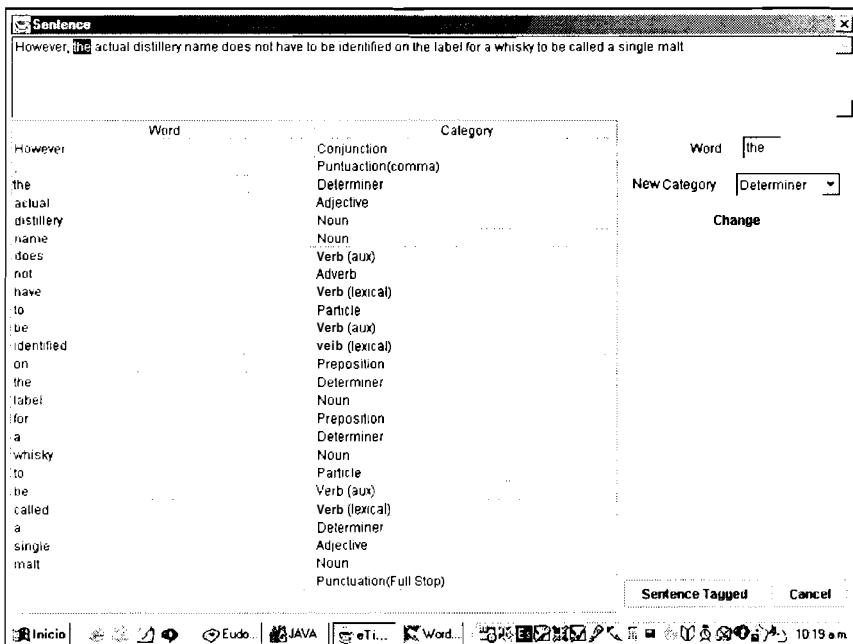


Figure 16. eTiKeT@ (Data Base Layout)

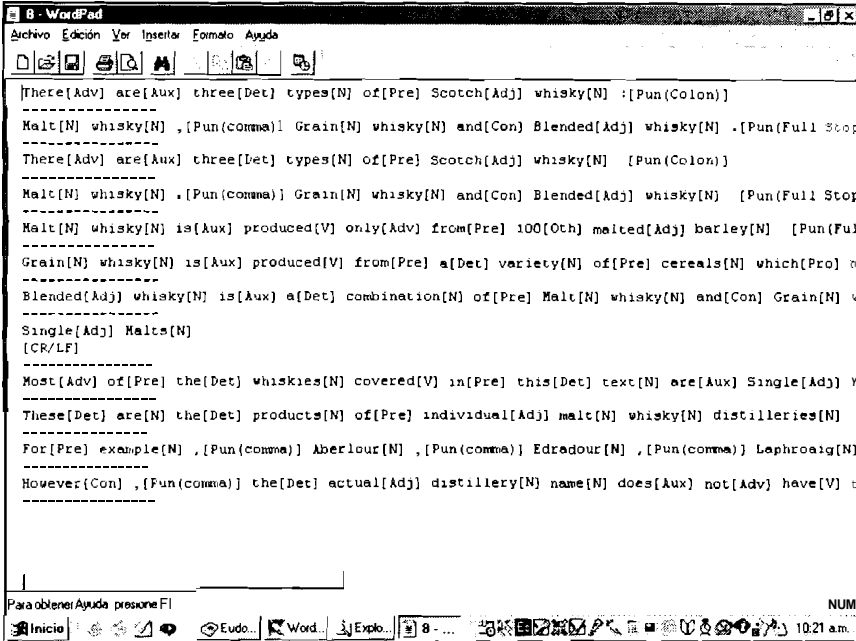


Figure 17. eTiKeT@ (ASCII Layout)

Microsoft Access - [Patrones: Tabla]

Archivo Edición Ver Insertar Formato Registros Herramientas Vexkana ?

PatronCode	Numante	Anteriores	Numpost	Posteriores	Cateyoria	Apariciones	SesionCode
1	1-2		32100	0	0	1	8
2	2-25		31006	2	1	1	8
3	3-252		3064	10	1	1	8
4	35210		3640	0	1	1	8
5	32100		340-2	6	1	1	8
6	31006		30-2-2	4	2	1	8
7	3064		2-2-2	0	2	1	8
8	1-2		30-10	0	1	1	8
9	2-20		3-100	0	1	1	8
10	300-1		3084	0	1	1	8
11	30-10		3840	0	1	1	8
12	3-100		340-2	8	1	1	8
13	3008		30-2-2	4	1	1	8
14	3084		2-2-2	0	1	1	8
15	1-2		3021	0	2	1	8
16	2-20		3215	0	1	1	8
17	3-200		3156	2	1	1	8
18	3002		35613	1	1	1	8
19	3021		36134	5	1	1	8
20	3215		31340	6	1	1	8
21	3156		340-2	13	1	1	8
22	35613		30-2-2	4	1	1	8
23	36134		2-2-2	0	1	1	8
24	2-20		3216	0	1	1	8
25	3-200		31610	2	1	1	8

Registro: 14 de 110

Vista Hoja de datos

NUM 10:24 a.m.

Figure 18. Information on Inferred Patterns and Statistics

The screenshot shows the TITag application window with a text document titled 'Novela.txt'. The text is a paragraph in Spanish. Below the text is a table with columns: 'Chk', 'Forma', 'Lema', 'Categoría', 'Flexión', and 'P. Nota'. The table lists various words from the text and their corresponding grammatical categories and inflections.

Chk	Forma	Lema	Categoría	Flexión	P. Nota
	EL	el	Artículo Determinado	MS	
	DESEO	deseo	Nombre Común	MS	
	DE	de	Preposición	Ø	
	DIOS	dios	Nombre Común	MS	
	Sus	su	Determinante Posesivo	CP	
	ojos	ojo	Nombre Común	MP	
	se	se	Pronombre Reflexivo	3ª PERS	
	oscurecien	oscurecer	Verbo Principal Normal	3ª P. IMP. IND	
	de	de	Preposición	Ø	
	manera	manera	Nombre Común	FS	
	intermitente	intermitente	Adjetivo Calificativo Normal	MS	
	#		Puntuación Punto y Seguido	Ø	
	#		Final Oración		
	Su	su	Determinante Posesivo	CS	
	mirado	mirada	Nombre Común	FS	
	languido	languido	Adjetivo Calificativo Normal	MS	
	contrastado	contrastar	Verbo Principal Normal	1ª S. IMP. IND	
	con	con	Preposición	Ø	
	la	el	Artículo Determinado	FS	
	lucidez	lucidez	Nombre Común	FS	
	de	de	Preposición	Ø	
	su	su	Determinante Posesivo	CS	
	mente	mente	Nombre Común	MS	
	#		Puntuación Punto y Seguido	Ø	
	#		Final Oración		
	El	el	Artículo Determinado	MS	

Figure 19. TL-Tag (POS-Tagger and Lemmatizer)

III.1.3. Lemmatisation

The distinction between words (tokens), the word forms (types) and base forms (lemmas) is important. Consider the following word sequence: *plays, playing, played, play, plays, play, playing, played* and *played*, where we have nine words (tokens), four word forms (types) and one lemma, namely *play*. As mentioned above, standard concordancers would process inflected forms (tokens) of the same base form (lemma) as different word forms (types). A way of dealing with this and other potential problems (see *POS-section*), which can seriously affect the counting of linguistic items, is to classify together all the identical or related forms of a word under a common headword: lemmatisation; just as in a dictionary where the various morphological inflected and derived forms of a word are listed under a single entry. In order to handle the complexities of morphology, including irregularities, lemmatisers typically employ two different but combined processes: (1) *εσϑε-βγ-οσϑε* method to deal with irregularities, by means of rules; for example, *better* and *best* are listed and counted under the headword *good*; (2) affix stripping method; if a word form is not listed under any headword (*οσϑε-βγ-οσϑε* method), then a number of affix stripping rules are applied; for example, the plural suffix *-s* is taken off the word form *cars*, outputting the base form *car*. Finally, if a word form does not appear as part of the affix rule

system, or is not listed as a specific exception, then the word is listed as a lemma in the same form in which it appears in the text. The lemmatisation process is normally performed automatically as part of the POS-process, producing enlarged tagged data lists: token, tag and lemma and grammatical information (gender, number, tense, etc.), i.e. *TL-Tag*¹² (Figure 19).

A useful CALL application based on the lemmatisation process is *Verbos Españoles Conjugados*¹³ (*VEC*). *VEC* has been designed to assist students in the correct use and spelling of Spanish verbs. It can be used as a stand alone program or run parallel as a grammatical help tool; the student just needs to write any Spanish verb form and *VEC* feeds back with full information on tense, mood, person and number (Figures 20 and 21).

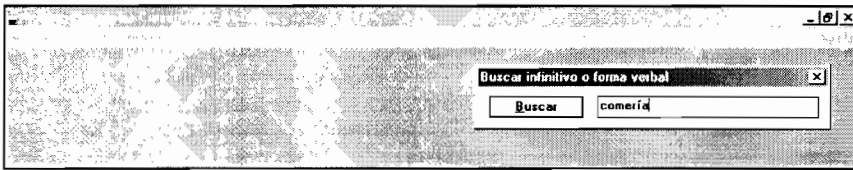


Figure 20. Student's Query

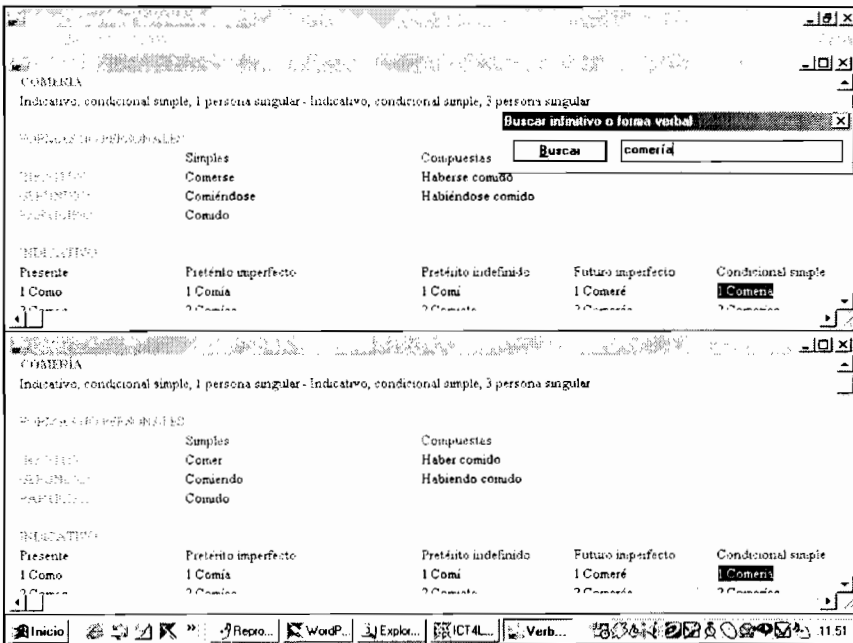


Figure 21. Full Verb Information

A much more interesting and challenging application is *Tagcorder*¹⁴. *Tagcorder* has been implemented to allow complex searches within the *CUMBRE Corpus* (Sánchez et al. 1995) and takes full advantage of tagged and lemmatised data. Users can invariably look for terminal nodes (types; *Figure 22*), non-terminal nodes (POS-tags; *Figure 23*) with any additional tagged grammatical information (*number, person, mood, tense, etc.*), base forms (lemmas; *Figure 24*) and/or any combinations of types, POS-tags and lemmas; *Figures 25 and 26*). The program itself is very interactive and flexible in its search procedure and extremely fast as it works with pre-indexed text.

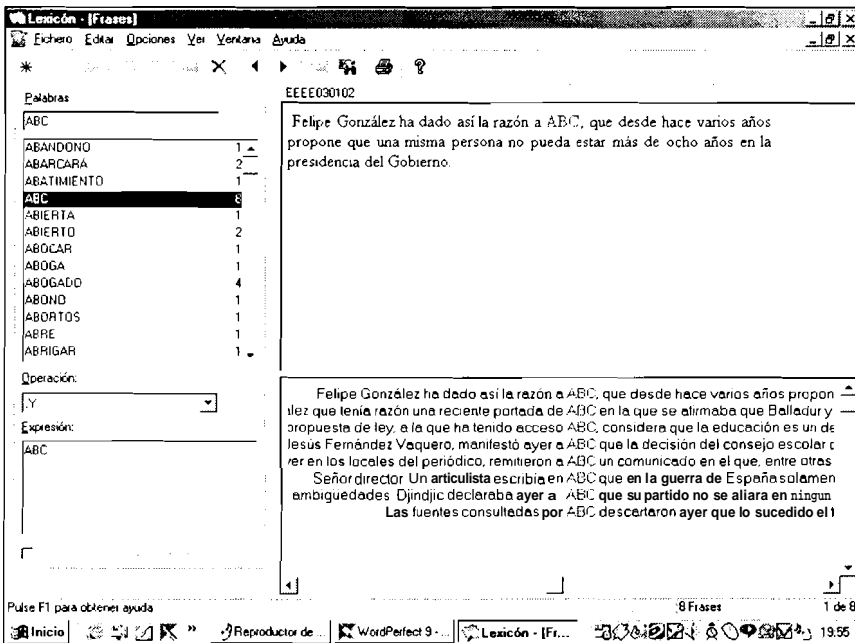


Figure 22. Type Search "ABC"

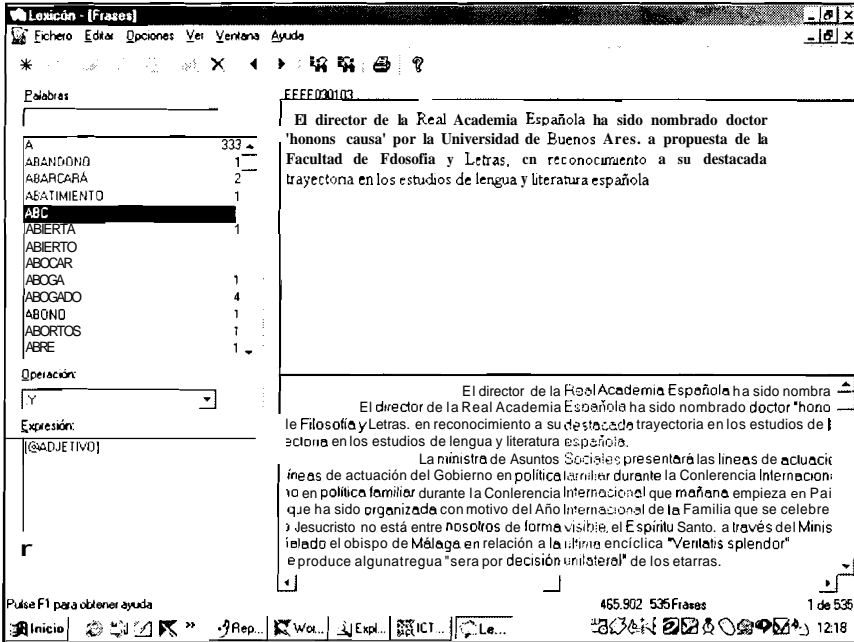


Figure 23. POS-Search "ADJECTIVE"

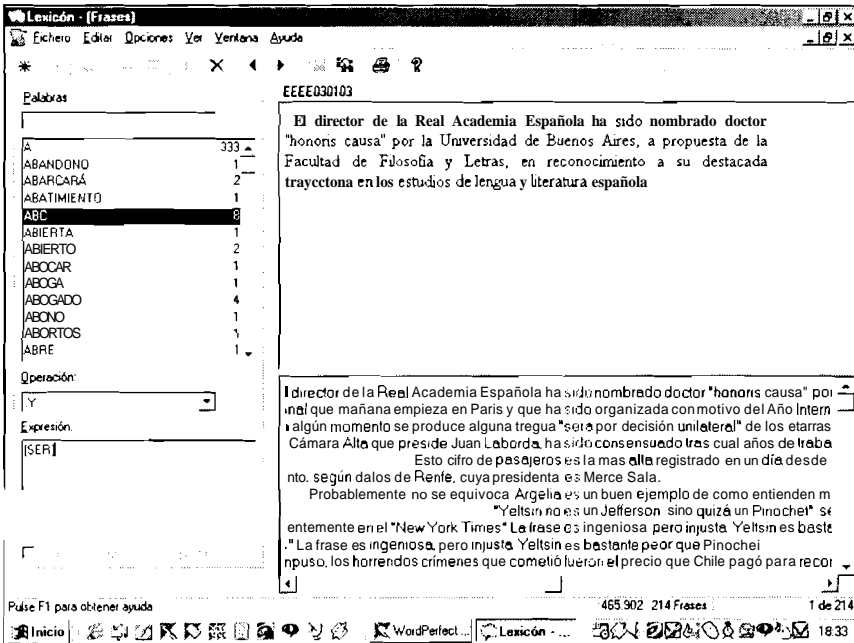


Figure 24. Lemina Search "SER"

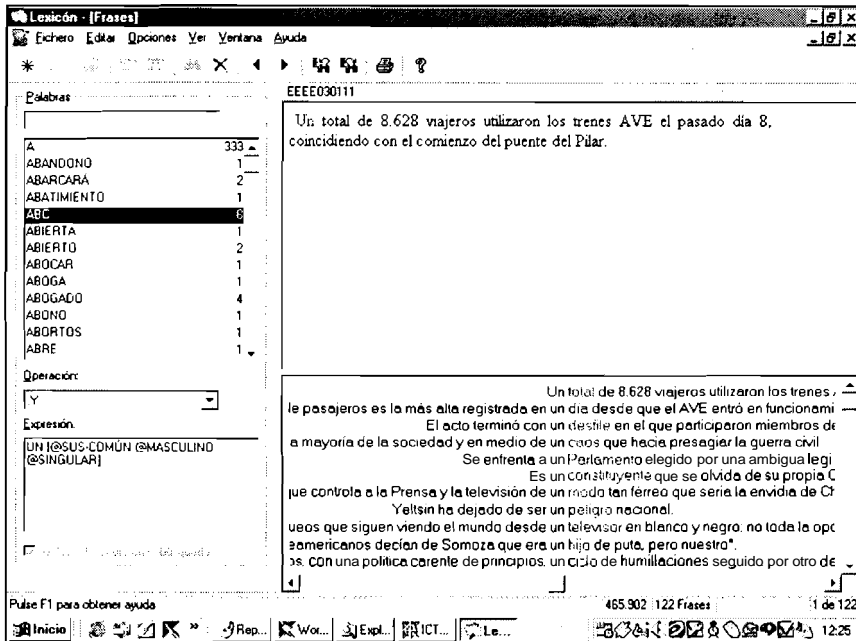


Figure 25. Complex Search: "UN" + NOUN (Countable + Masculine + Singular)

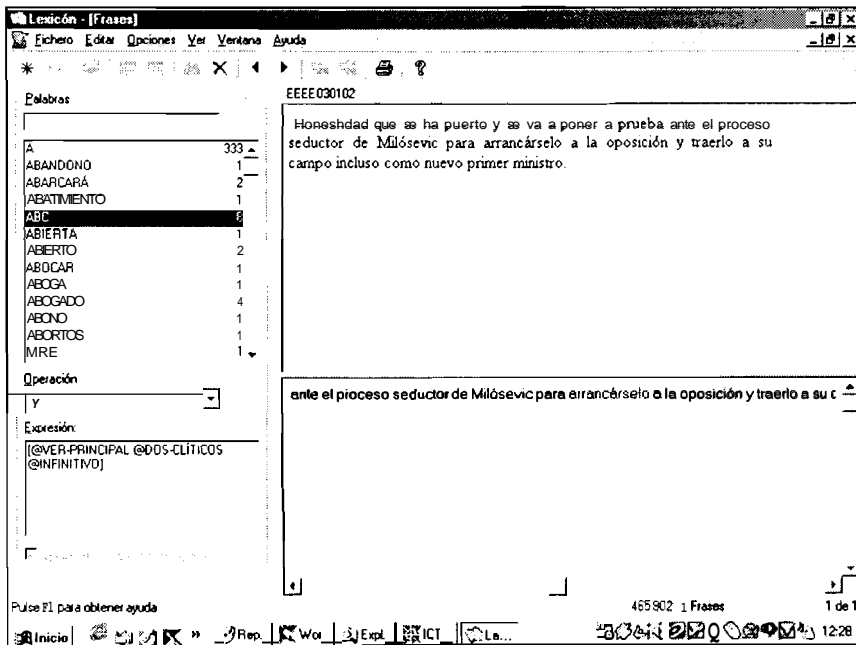


Figure 26. Complex Search: VERB (Main + Two Clitics + Infinitive)

III.1.4. Parsing

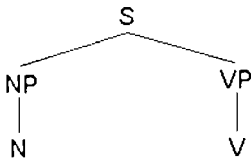
Parsing involves the procedure of bringing basic morphosyntactic categories into high-level syntactic relationships with one another. This is probably the most commonly encountered form of corpus annotation after POS tagging. Parsed corpora are sometimes known as *treebanks*.

There are rules governing the way in which words can be put together to form syntactically well-formed or grammatical sentences: the study of syntax aims to discover them and to describe and analyse language in terms of these rules. Consider the sentence *A dog chased that girl*, where we find the same pattern of constituents before and after the verb. *that* is *determiner + noun*. These two words also appear to belong together more closely than say the noun *dog* and the verb *chased*. Another way of illustrating that these words belong together is to give the *girl* and the *dog* a name—names of specific items such as individual people, animals, places and so on called *proper nouns*—, and we get, for example *Henry chased Carol*.

It seems clear that natural languages or human languages have a role of *constituent structure*. A sentence is not just a mere string of words. The words are grouped into phrases, each of which consists of a short phrase. Many of the important properties of languages are organised around constituent structure. Constituent structures (a) group words into constituents such as *the rlog* and *into the garden*; (b) give names to the constituents, such as *noun phrase* and *prepositional phrase*. In turn, constituent structures are sanctioned or generated by rules, known as *phrase-structure rules* of this type:

$$\begin{aligned} S &\rightarrow NP VP \\ NP &\rightarrow N \\ VT &\rightarrow V \end{aligned}$$

where *S* stands for sentence, *NP* for noun phrase, *VP* for verb phrase, *N* for noun and *V* for verb. So the PS-rules above state that a sentence consists of a noun phrase followed by a verb phrase. In turn, the *NP* of an *N* and the *VP* of a single *V*. The tree structure derived or generated by that rule would be



Parsing algorithms can proceed top-down or bottom-up. In some cases, top-down and bottom-up algorithms can be combined¹⁵.

The Visual Interactive Syntax Learning (VISL) website¹⁶ is particularly interesting and useful for language learners. It contains an on-line parser and a variety of other tools concerned with English grammar, including games and quizzes. The parser itself is an excellent and very transparent application that allows learners to analyse and experiment on sentences and study their structure (Figure 27).

Interesting in this respect is also the parsing of students' erroneous input. Integrated parsers into CALL software can be prepared to deal with linguistic errors in the input. So the grammar that copes with correct sentences is complemented with a grammar of incorrect sentences. The advantage of this error grammar approach is that the feedback to students' output can be very specific and is normally fairly reliable as it can be attached to a very specific rule. However, the major drawback of this approach is that individual learner errors have to be anticipated in the sense that each error needs to be covered by an adequate rule.

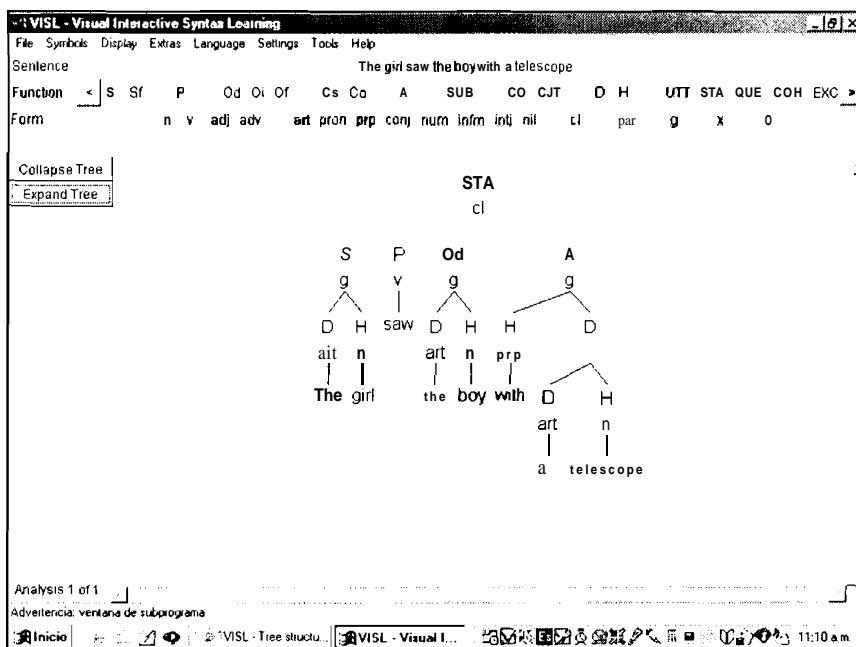


Figure 27. Visual Interactive On-Line Parser

III.1.5. Speech Technology¹⁷

CALL software has normally been restricted to written text. However, recent advances in multimedia have resulted into powerful hardware and software applications, allowing users to

attach a microphone and loudspeakers to soundcards and to record his/her own voice. Furthermore, storing these sound files is no more a problem due to the immensely increased capacity and cost reduction of hard disks, other storage devices (CD-ROM) and improved compression algorithms for this kind of data (i.e. MP3 files).

Presently, there is a wide range of speech software available. This includes (a) spoken input processing¹⁸ or speech analysis, where speech input is analysed and represented graphically or numerically; (2) speech recognition: the transformation of spoken input into written output; and (3) speech synthesis, that is the conversion of text to speech¹⁹; this includes not just matching characters to sounds, but also intonation and the rhythm particular utterances have. Advances in speech synthesis technology have reached a high level of performance and robustness and some CALL applications have started considering its integration²⁰.

In contrast, speech recognition is far more complex than speech synthesis. Speech recognition needs an extensive analysis of speech by means of a number of parameters, which are very difficult to establish as they can be easily affected by background noise, speech speed (connected speech), particular accents or idiosyncratic individual's speech. All this leads to complicate and interfere in the fixing and interpretation of the established parameters.

There are some commercial applications able to "understand" natural speech and can provide language students with realistic, highly effective, and motivating speech practice. One of this application is *IBM® ViaVoice®*. This program runs on normal PCs and includes speaker independent continuous speech recognition engines and is able to deal with complete sentences spoken at a natural pace, not just isolated words, though it requires a minor training period. To run the program, the user just needs to associate it with any wordprocessor, where the user's utterances will to be transcribed in (Figure 28) or run *Speech Pad*, a simplified standard wordprocessor that includes *ViaVoice*.

Many multimedia CALL courses already have and still include some naive direct pronunciation practice. That is, exercises which focus on pronunciation, fluency and word order, and with native speaker models which are heard immediately after a student's performance. These applications leave the learner-model comparison to student's criteria or visualise graphically both performances, indicating the success rate in %. The negative side of these exercises is that in some instances it is even for naive speakers of the language very difficult, if not impossible, to achieve satisfactory success rates. Some of these applications are neither very flexible nor accurate and sometimes students would need to repeat their utterances in several occasions before the program "understands" them correctly. In turn, this might lead to some small frustration.

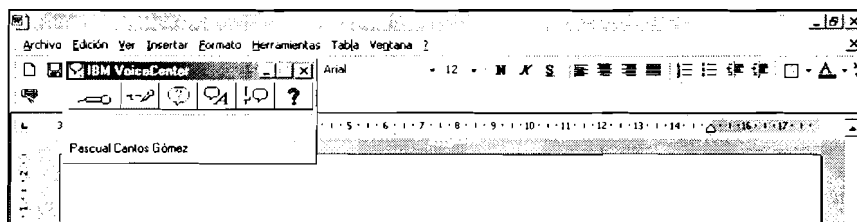


Figure 28. IBM ViaVoice

V. CONCLUSION

CALL has for long been dominated by the drills and practice associated with behaviourism (cloze and gap-filling exercises, multiple-choice tests, etc.) and, eventually, the use of some basic word-processing tools. Soon, some language teachers and CALL practitioners reacted negatively and noted a lack of progress in CALL (Kaliski 1993, Last 1992, etc.), partially due to the:

- Limited CALL software available
- Reduced number of computational exercises
- Incompatibility between employed CALL techniques and current language teaching pedagogy (particularly influential in this respect has been the emergence of communicative syllabus)
- Consequence of new technology being unable to fulfill teachers' expectations

Fortunately, things have changed for CALL in the 90s, partly because of the wider availability of PCs and the integration of linguistic corpora and NLP-technology. The use of linguistic corpora and NLP-applications are highly valuable tools for language description with important implications for language teaching, as they can:

- Assist language teacher in identifying relevant content of instruction (vocabulary, grammar, contexts, etc.), and
- Help in developing new pedagogical and methodological approaches to instruction (i.e. shifting the pedagogical teaching/learning paradigm from computer as manager to computer as pedagogue).

Modern CALL has changed and instead of adapting it to what software can offer, an attempt is made to get it to take account of the necessary conditions of successful language learning. Learners are given much more control over what they learn: autonomous language learning is self-paced, more interactive, meaning dominated, task-oriented activities (Kennedy

1998: 393).

Particularly interesting in this respect is the use of real and relevant text sample for students and teachers as the central pedagogical teaching/learning cornerstone. Real time manipulation of texts by students using integrated user-friendly interfaces, including word-processing tools and NLP-applications could conform an extremely valuable pedagogical paradigm within the foreign language learning/teaching context. Teachers would be able:

- To extract, manipulate and adapt texts to students needs and language level
- To enrich plain texts with POS-tags and syntactic annotations for class work
- To extract vocabulary lists, phrase lists, concordance lists, etc. (sublanguage specific, adapted to a specific level or domain, etc.)
- To generate automatically *ad hoc* exercises, depending on students' particular needs.

Similarly, students could also take advantage of this integrative CALL application

- To explore the target language by means of concordancers with integrated taggers and parsers and/or tagged texts and treebanks
- To extract the gist of more difficult texts, using MT-software
- To check the meaning of words and phrases (electronic dictionaries)
- To generate automatically *ad hoc* exercise generation, depending on one's own needs
- To hear the text, selected sentences or words, using speech synthesizers
- To answer orally to some responses, dictating the solutions to the computer (speech recognition tools)

Actually, what we propose is a sophisticated CALL language processing tool²¹ that

- Takes full advantage of current computational advances in an integrated and unitary way:
 - Electronic dictionaries (monolingual, bilingual or multilingual ones)
 - MT systems
 - POS-taggers
 - Syntactic parsers
 - Concordancers
 - Speech production/recognition
 - Word-processors (this includes spell-checkers and grammar and style checkers)
- Goes beyond written text, as it also accounts for oral production and oral recognition
- Assists both teachers and students in their respective tasks and that could contribute to new and challenging pedagogical and methodological paradigms in the area of foreign language learning.

And since we have all these computational tools at our disposal, it makes no sense to renounce their application in such an important area as language pedagogy. We cannot dismiss them, we must use them ...

NOTES

1. See among others, the Oxford Concordance Program (OUP 1988), Longman Mini Concordancer (Chandler and Tribble 1989), MicroConcord (Scott and Johns 1993), MonoConc (Barlow 1996), TACT (Bradley and Presutti 1989), or WordSmith (Scott 1999).
2. <http://web.bham.ac.uk/~johnstf/timeone.htm>
3. *Context* can be downloaded free of charge for non-commercial purposes from Tim Johns DDL-web page: <http://web.bham.ac.uk/~johnstf/timeone.htm>
4. Sánchez, A. and P. Caiitos (3000) *Practica tu vocabulario*. Madrid: SCEL.
5. For a more comprehensive survey, visit "Module 3.5. Human Language Technologies (HLT)" of the *Information and Communications Technology for Language Teacher (ICT4LT)* web page: http://www.ict4lt.org/en/en_mod3-5.htm
6. Hutchins and Somers (1992) and Arnold et al. (1994) provide excellent introductions to MT.
7. The MT system used here is *Spanish Assistant* (MicroTac Software). Other commercial MT software: *Systran* (<http://babelfish.altavista.digital.com/>) or *Power Translator* (<http://www.lhsl.com/powertranslator/>).
8. See Greene and Rubin (1971) for a detailed description of this tagger.
9. Described in detail by Garside (1987) and Marshall (1987).
10. To get a free copy, for academic purposes only, e-mail Rafael Valencia (rafavalencia@ono.es), Rodrigo Martínez (rodrigo@dif.um.es) or Pascual Caiitos (pcantos@um.es).
11. Where 0 = English and 1 = Spanish.
13. *TL-Tag* (TechnoLingua) is part of the *CUMBRE Corpus* Project and is not yet commercially available. For those interested in it contact any of the people involved in the Project: Enrique Pérez de Lema (delema@jazzfree.com), José Simón (jsj38746@telefonos.es), Aquilino Sáizcliez (asanchez@um.es) or Pascual Cantos (pcantos@um.es).
13. Diez, P. L. and J. Iborra (1909) *Verbos Españoles Conjugados*. Madrid: SGEL.
14. See also note 12.
15. Allen (1995), Conington (1994) and Gazdar and Mellish (1989) include excellent introductory sections on parsing algorithms.
16. <http://visl.hum.ou.dk/>
17. An excellent site is Integrating Speech Technology in (Language) Learning: <http://www.instil.org>.
18. Among the most interesting web sites, check: <http://agoralang.com/signalvze.html>.
19. A good example is *Winspeech*: <http://www.pcww.com>.
20. The Polytechnic University of Hong Kong site includes a number of text-to-speech tools: <http://vlc.polyu.edu.hk/TextToSpeech>
21. We have deliberately not considered the integration of *Information Technologies* here. This would have inevitably expanded the potential of the "tool" proposed here.

REFERENCES

- Allen, J. (1995). *Natural Language Understanding*. Redwood City, CA: The Benjamin/Cumming.
- Arnold, D. et al. (1994). *Machine Translation*. Oxford: NCC Blackwell.
- Aston, G. (1995). Corpora in Language Pedagogy: Matching Theory and Practice. In G. Cook and B. Seidlhofer (Eds.), *Principle and Practice in Applied Linguistics: Studies in Honour of H.G. Widdowson*. Oxford: Oxford University Press, 257-270.
- Aston, G. (1996). The British National Corpus as a Language Learner Resource. In S. Botley, J. Glass, T. McEnery and A. Wilson (Eds.), *Proceedings of Teaching and Language Corpora 1996*. Lancaster: Lancaster University Press, 178-191.
- Aston, G. (1997). Small and Large Corpora in Language Learning. In B. Lewandowska-Toinaszczyk and J. Melia (Eds.), *Proceedings of the First International Conference on Practical Applications in Language Corpora*. Lodz: Lodz University Press, 51-62.
- Ball, C. (1995). Hwæt! A Corpus-based Hypermedia Resource for Old English Lexical Acquisition. *Linguistic Society of America Annual Meeting*, New Orleans, January 7, 1995 (Software poster session).
- Barlow, M. (1992). Using Concordance Software in Language Teaching and Research. *Second International Conference on Foreign Language Education and Technology*. Kasugai, Japan (Conference).
- Barlow, M. (1995). Corpora for Theory and Practice. *International Journal of Corpus Linguistics*, 1(1), 1-38.
- Bazerman, C. (1994). *Constructing Experience*. Carbondale: Southern Illinois University Press.
- Brown P.F. et al. (1993). The Mathematics of Statistical Machine Translation: Parameter Estimation. *Computational Linguistics*, 19(2), 263-311.
- Burnard, L. (1995). *The British National Corpus Users Reference Guide*. Oxford: Oxford University Computing Services.
- Burnard, L. and T. McEnery (Eds.) (1999). *Rethinking Language Pedagogy from a Corpus Perspective*. Hamburg: Peter Lang.
- Celce-Murcia, M. (1990). Data-based Language Analysis and TESL. In J. E. Alatis (Ed.), *Georgetown University Round Table on Language and Linguistics 1990*. Georgetown: Georgetown University Press.
- Cobb, T. (1997). Is there any Measurable Learning from Hints-only Concordancing?. *System*, 25(3), 301-315.
- Cole, R. (1996). Foreword. In R. Cole (Ed.), *Survey of the State of the Art in Human Language Technology*. In <http://eslu.cse.ogi.edu/HLTsurvey/HLTsurvey.html>.

- Collins, H. (1999). Materials Design and Language Corpora: a Report in the Context of Distance Education. In L. Burnard and T. McEnery (Eds.), *Rethinking Language Pedagogy from a Corpus Perspective*. Hamburg: Peter Lang, 51-63.
- Covington, M. (1994). *Natural Language Processing for Prolog Programmers*. Englewood Cliffs, NJ: Prentice Hall.
- Flowerdew, J. (1993). Concordancing as a Tool in Course Design. *System* 21(3), 231-243
- Flowerdew, J. (1996). Concordancing in Language Learning. In M. Pennington (Ed.) *The Power of CALL*. Houston: Athelstan, 97-113.
- Garside, R. (1987). The CLAWS Word-tagging System. In R. Garside et al. (Eds.), *The Computational Analysis of English. A Corpus-Based Approach*. London: Longman, 30-41.
- Gavioli, L. (1997). Exploring Tests through the Concordancer: Guiding the Learner. In A. Wichmann, S. Fligelstone, T. McEnery and G. Knowles (Eds.), *Teaching and Language Corpora*. London: Longman, 83-99.
- Gazdar, G. and C. Mellish (1989). *Natural Language Processing in PROLOG. An Introduction to Computational Linguistics*. Wokingham: Addison-Wesley.
- Greene, B. and G. M. Rubin (1971). *Automated Grammatical Tagging of English*. Providence Rhode Island: Brown University Press.
- Higgins, J. (1988). *Language, Learners and Computers*. London: Longman.
- Higgins, J. (1991a). Looking for Patterns. In T. Jollis and P. King (Eds.), *Classroom Concordancing*. Birmingham: Birmingham University Press, 63-70.
- Higgins, J. (1991b). Which Concordancer? A Comparative Review of MS-DOS Software. *System*, 19(1/2), 91-100.
- Hutchins, J. and H. Somers (1992). *An Introduction to Machine Translation*. London: Harcourt Brace Jovanovich Publishers.
- Johansson, S. (1980). The LOB Corpus of British English Texts: Presentation and Comments. *ALLC Journal*, 1, 25-36.
- Jollis, T. (1986). Microconcord: a Language-learner's Research Tool. *System*, 14(2), 151-162
- Jollis, T. (1988). Whence and Whither Classroom Concordancing?. In T. Bongaerts, P. de Haan, S. Lobbe, and H. Wekker (Eds.), *Computer Applications in Language Learning*. Dordrecht: Foris, 9-27.
- Jollis, T. (1991). Should you be Persuaded—Two Samples of Data-driven Learning Materials. In T. Jollis and P. King (Eds.), *Classroom Concordancing*. Birmingham: Birmingham University Press, 1-13.

- Johns, T. (1993). Data-driven Learning: an Update. *TELL&CALL*, 1993/2., 4-10.
- Johns, T and P. King (Eds.)(1991). *Classroom Concordancing*. Birmingham: Birmingham University Press.
- Kacowsky, W. (1987). *Better English*. Salzburg: Jugend-Verlag.
- Kaliski, T. (1992). Computer-assisted Language Learning. In P. Roach (Ed), *Computing in Linguistics and Phonetics: Introductory Readings*. London: Academic Press, 97-106.
- Kennedy, G. (1998). *An Introduction to Corpus Linguistics*. London: Longman.
- Kettemann, B. (1995). On the Use of Concordancing in ELT. *TELL&CALL*, 1995/4, 4-15
- Kučera, H. and N. Francis (1967). *Computational Analysis of Present Day American English*. Providence Rhode Island: Brown University Press.
- Last, R. (1992). Computers and Language Learning: Past, Present — and the Future?. In C. Butler (Ed.). *Computers and Written Texts*. Oxford: Blackwells, 227-245.
- Leech, G and C. N. Candlin (Eds.)(1986). *Computers in English Language Teaching and Research*. London: Longman.
- Marshall, I. (1987). Tag Selection Using Probabilistic Methods. In R. Garside et al. (Eds.), *The Computational Analysis of English. A Corpus-Based Approach*. London: Longman, 42-56.
- Nagao, M. (1984). A Framework of a Mechanical Translation between Japanese and English by Analogy Principle. In A. Elithorn and R. Banerji (Eds.). *Artificial and Human Intelligence*. Amsterdam: Elsevier Science Publisher, 173-180.
- Sánchez, A. et al. (1995). *CUMBRE corpus lingüístico del español contemporáneo. Fundamentos, metodología y aplicaciones*. Madrid: SGEL.
- Sinclair, J. (1991). *Corpus, Concordance and Collocation*. Oxford: Oxford University Press
- Stevens, V. (2005). Concordancing with Language Learners: Why? When? What?. *CAELL Journal*, 6(2), 2-10.
- Tribble, C. (2007). Improving Corpora for ELT: Quick and Dirty Ways of Developing Corpora for Language Teaching. In B. Lewandowska-Tomaszczyk and J. Melia (Eds.), *Proceedings of the First International Conference on Practical Applications in Language Corpora*. Lodz: Lodz University Press. 106-117.
- Tribble, C. and G. Jones (1990). *Concordances in the Classroom*. London: Longman.
- Widdowson, H.R. (1983). New Starts and Different Kinds of Failure. In D. Larsen-Freeman et al. (Eds.), *Learning to Write: First Language/Second Language*, New York, Longman, 34-47.