

# Optimizing the role of language in Technology-Enhanced Learning

Sylviane Granger

► **To cite this version:**

Sylviane Granger. Optimizing the role of language in Technology-Enhanced Learning. Granger, Sylviane. NOE-Kaleidoscope, 71 p., 2007. hal-00197203

**HAL Id: hal-00197203**

**<https://telearn.archives-ouvertes.fr/hal-00197203>**

Submitted on 14 Dec 2007

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



**kaleidoscope**  
shaping the scientific evolution of Technology Enhanced Learning

# **Optimizing the role of language in Technology-Enhanced Learning**

Sylviane Granger (ed.)

Louvain-la-Neuve, 4-5 October 2007

## **SCIENTIFIC COMMITTEE**

Georges Antoniadis (Université Stendhal Grenoble, France)  
Cédric Fairon (Université catholique de Louvain, Belgium)  
Sylviane Granger (Université catholique de Louvain, Belgium)  
Adam Kilgarriff (Lexical Computing Ltd, United Kingdom)  
Fanny Meunier (Université catholique de Louvain, Belgium)  
John Nerbonne (Universiteit Groningen, The Netherlands)  
Magali Paquot (Université catholique de Louvain, Belgium)  
Claude Ponton (Université Stendhal Grenoble, France)  
Michael Rundell (Lexicography Masterclass Ltd & Macmillan Dictionaries, United Kingdom)  
Serge Verlinde (Katholieke Universiteit Leuven, Belgium)

## **ORGANIZING COMMITTEE**

Centre for English Corpus Linguistics, Université catholique de Louvain (Belgium)

Sylviane Granger  
Fanny Meunier  
Sylvie De Cock  
Valérie Decuyper  
Céline Gouverneur  
Claire Hugon  
Marie-Aude Lefer  
Magali Paquot  
Jennifer Thewissen

Centre de traitement automatique du langage, Université catholique de Louvain (Belgium)

Cédric Fairon  
Bernadette Dehottay  
Thomas François



## Table of Contents

|  |    |
|--|----|
| Charles Alderson, <i>Computer-adaptive language testing</i> .....  | 1  |
| Eckhard Bick, <i>VISL: A cross-language approach to NLP- and games-based grammar teaching</i> .....                                      | 5  |
| Jozef Colpaert, <i>CALL software design principles and the integration of NLP</i> .....  | 7  |
| Piet Desmet and Hans Paulussen, <i>CorpusCALL: Challenges and opportunities</i> ...  | 9  |
| Sylviane Granger, <i>The contribution of learner corpus research to TELL</i> .....   | 13 |
| Thomas Hansen, <i>Feedback methods in computer assisted pronunciation training applications using automatic speech recognition</i> ..... | 17 |
| Holger Hvelplund, <i>Language technology projects at IDM</i> .....   | 19 |
| Adam Kilgarriff, <i>Using corpora in language learning: the Sketch Engine</i> .....  | 21 |
| Thomas Koller, <i>A plurilingual ICALL system for Romance languages</i> .....  | 25 |
| Agnes Kukulka-Hulme, <i>Like stars in the firmament: Language learning on mobile devices</i> .....                                       | 29 |
| Marie-Noëlle Lamy, <i>Computer-mediated communication for language learning</i> .....  | 33 |
| Claudia Leacock, <i>Writing English as a second language: A proofreading tool</i> ....   | 37 |
| Fanny Meunier, <i>Second language acquisition theory and TELL</i> .....  | 41 |
| John Nerbonne, <i>Detecting syntactic interference</i> .....   | 43 |
| Richard Pemberton, <i>Planning a smart phone system to support self-directed L2 vocabulary learning</i> .....                            | 45 |
| Michael Rundell, <i>The dictionary of the future</i> .....   | 49 |
| Emma Shercliff, <i>Macmillan English Campus: A case study</i> .....  | 53 |
| Serge Verlinde, <i>The Base lexicale du français (BLF), a free Web-based learning environment for French vocabulary</i> .....            | 57 |
| Carl Vogel, <i>Linguistic anomaly</i> .....  | 61 |
| David Wible, <i>Evolving approaches to Web-supported language learning: From platforms to platform-independent tools</i> .....           | 65 |
| Johannes Widmann, <i>The SACODEYL project – Corpus exploitation for language learning purposes</i> .....                                 | 69 |



# Computer-adaptive language testing

Charles Alderson  
Lancaster University

## 1. Introduction

Computer-adaptive testing (CAT) involves presenting learners with items thought to be most suitable for them, and adjusting the selection of items in light of the learners' responses to previous items. The classic case of CAT involves the construction of a bank (collection) of test items which have been calibrated in terms of their empirical difficulty. Learners are typically initially presented with an item of medium difficulty. If their response is correct, they will then be presented with a more difficult item. If their response is incorrect, they are then presented with an easier item. If their response to the second item is correct they are given a more difficult item and if incorrect, an easier item. The computer calculates the learner's ability level (or score) on the fly as well as the reliability of the test as administered up to that point. Items from the bank are presented to test-takers following specially developed algorithms for the selection of the initial test item, subsequent test items, and a rule for concluding the test – ie the criteria to be met for the test to be terminated. Typically, the test is terminated when a given level of reliability has been reached, or when a pre-determined number of items has been delivered.

The advantages are that tests can be tailored to a learner's ability level rather than wasting time and effort by presenting them with items that are far too easy or far too difficult. As a consequence, tests can be markedly shorter than traditional linear tests, and thus more efficient. In addition, since each learner takes a different test than his or her fellow test-takers, cheating is made much more difficult. The major disadvantages are that in order to be able to predict a learner's ability level items need to be pre-tested and analysed using an Item Response Theory model – IRT (which allows the estimation of a learner's ability level independent of the difficulty of the items) – but IRT requires relatively large numbers of pilot test candidates for reliable ability estimates. Secondly, in high-stakes testing situations (like the TOEFL) learners are often schooled in remembering which items they have taken, and the item bank can be reconstructed if sufficient numbers of candidates recall the items (this has happened in China, for example, where CAT versions of TOEFL were compromised). In addition, truth-in-testing laws in some states of the USA mean that the test items constituting the basis of the test-taker's score have to be made available to test-takers on request.

This inevitably compromises the test bank. As a result, ETS (the developers of TOEFL) have ceased developing CATS since, as they put it, “feeding the CAT is too expensive”.

## **2. Outline of the presentation/demo**

In this talk I will briefly present the work I have done on CATs to date and then discuss possible amendments to the design of CATs to make them more relevant to learners and to test purposes. My main involvement with CATs, apart from critiquing their value, has been on the DIALANG Project. DIALANG is a suite of Internet-delivered diagnostic tests of 14 European languages, funded by the European Union, and based on the Common European Framework. It contains tests of reading, listening, writing, grammar and vocabulary, as well as a test of vocabulary size and a self assessment battery. The current version of the test is adaptive at what we call the test level, and I shall describe how this works and how feedback is given. Algorithms have been developed – but not yet implemented – to make the test adaptive at the item level, as in the classic case described above, and I shall describe the problems and solutions involved in making a test, which is delivered over the Internet, adaptive at item level. I shall then describe and discuss further developments of CATs. These are much less likely to take place in the context of high-stakes proficiency tests, for the reasons given above, but in the context of language learning, be that of progress or achievement tests or, in the case I am interested in, in further exploration and refinement of the diagnosis of learners’ strengths and weaknesses.

Adaptation to the learner’s response need not be simply on the basis of the difficulty of the item responded to, but on the basis of item content. Thus, if the items in the CAT bank are characterised not merely in terms of their empirical difficulty, but in terms of the language features they test, or in terms of the skill or sub-skill they measure, then one can envisage an adjustment of the CAT algorithm to take account of what is being tested. Thus, a diagnostic CAT of one’s command of structures in the language could select items on, for example, the use of the present perfect, and explore how thoroughly a learner mastered that tense/ aspect in a variety of contexts, or with a range of different verbs. Success on a range of items could lead to the selection of items on a different aspect of syntax, or to the presentation of items in the same syntactic area but known to be more “advanced” in terms of the acquisition sequence and/or the development of the learner’s syntactic competence. Similarly, one could envisage tests of vocabulary being structured according to the frequency of words in the language, or according to particular semantic fields, or domains of use or register, and so on. Computer adaptivity would then enable a more or less thorough exploration of strengths and weaknesses in lexical knowledge.

Similarly, in tests designed to establish a learner’s level on the Common European Framework of Reference or some similar relevant standard, CATS could present items calibrated and standard-set at the different CEFR levels, and the degree of a learner’s



mastery of items in a given skill or language use domain at each particular level could be explored in some depth.

Another adaptation of the principle of adaptivity could take account of learner characteristics (age, mother tongue, years of learning, gender, topics of interest, area of academic study, etc). The learner would select from a menu of possible characteristics those that applied to them, and the computer would only present items known to be suitable for learners with such a profile – or, indeed, items known to be a challenge for such learners.

Finally, instead of the computer making the decision on which next item to select, the learner could be allowed to do so (by, for example, requesting a more difficult item, or one on a different linguistic feature or another topic or academic discipline).

### **3. Issues and challenges**

The major challenge in the field is to identify relevant characteristics of items and of learners which would provide meaningful diagnoses, or results relevant to further learning, and this requires a much better theory of diagnosis and language development than we currently possess. I hope in the discussion that we can explore whether NLP techniques can contribute to this notion of adaptivity.

### **4. References**

- CHALHOUB-DEVILLE M. (ed) (2000), *Computer-adaptive tests of reading*, Cambridge, CUP.
- CHALHOUB-DEVILLE M. and DEVILLE C. (1999), “Computer Adaptive Testing in Second Language Contexts” in *Annual Review of Applied Linguistics*, 19: 273-299.
- DUNKEL P.A. (1999), *Considerations in Developing and Using Computer-Adaptive Tests to Assess Second Language Proficiency*, <http://www.cal.org/resources/Digest/cat.html> (last accessed 7.8.07).
- WAINER H. (ed.) (2000), *Computer-Adaptive Testing: A Primer*, Mahwah, NJ, USA, Lawrence Erlbaum Associates.



# **VISL: A cross-language approach to NLP- and games-based grammar teaching**

Eckhard Bick  
University of Southern Denmark

## **1. Introduction**

VISL (Visual Interactive Syntax Learning) is an integrated interactive user interface for teaching grammatical analysis on the Internet, developed at the University of Southern Denmark, offering a unified system of analysis for 25 different languages, 8 of which are supported by live grammatical analysis of running text. For reasons of robustness, efficiency and correctness, the system's internal tools are based on the Constraint Grammar formalism (Karlsson 1990), but users are free to choose from a variety of notational filters, supporting different descriptive paradigms, with a current teaching focus on syntactic tree structures, language independent grammatical categories and the form-function dichotomy. VISL's core NLP-programs use the author's hybrid multi-level parsers (<http://beta.visl.sdu.dk>), while teaching applications (<http://visl.sdu.dk>) and corpus searching tools (<http://corp.hum.sdu.dk>) are implemented as platform independent Java-programs and Perl-cgi's. Though lexica and parsing rules are developed individually for each language, a common CG and treebank data format facilitates source data transfer into grammar teaching games, structural or color based visualisation, and linguistic revision of corpus data.

## **2. Outline of the presentation/demo**

In a modern school or university environment, grammar teaching is often plagued by the fact that the subject is perceived as "academic" and "uninteresting", and its inherent analytical view on language conflicts with a current language teaching focus on assimilation, naturalness, communicative media etc. Also, grammar teaching is affected by a cross-language handicap, because students are confronted with different formal systems and terminology, depending on the individual language taught (e.g. latinid vs. native terminology, morphological vs. functional word classes, syntax trees for English, dependency grammar for Czech, topological fields for Danish). In order to address these problems, the VISL system has introduced a novel, unified approach

across languages, built on a clear distinction between function and form, and tied to visually stable clues, such as iconic abbreviations, symbols and colour coding.

The presentation/demo will demonstrate how these principles can be implemented in the form of internet-based grammar games such as *WordFall*, *Labyrinth*, *Syntris* etc., as well as tree structures and corpus tools.

### 3. Issues and challenges

While games and treebanks can be based on manually annotated data, a truly flexible system and, not least, language teaching based on empirical, corpus-derived evidence cannot realise its full potential without robust, automatic NLP, and even apparently “closed” exercises and games become dependent on such tools if a higher degree of lexical or structural variation is to be achieved, or where a teacher would like to adapt exercises or games to a given text book tradition or recently treated literature. For these reasons (and also for the sake of linguistically more robust interfaces), I believe the integration of main stream parsing technology to be one of the major challenges in the future development of CALL applications.

### 4. References

- BICK E. (2005-1), “Grammar for Fun: IT-based Grammar Learning with VISL”, in P.J. Henriksen (ed.), *CALL for the Nordic Languages*, København, Samfundslitteratur: 49-64 (Copenhagen Studies in Language).
- BICK E. (2005-2), “Live use of Corpus data and Corpus annotation tools in CALL: Some new developments in VISL”, in H. Holmboe (red.), *Nordic Language Technology, Årbog for Nordisk Sprogteknologisk Forskningsprogram 2000-2004 (Yearbook 2004)*, Copenhagen, Museum Tusulanum: 171-186.
- DAVIES G. (ed.) (2007), *Information and Communications Technology for Language Teachers (ICT4LT)*, Slough, Thames Valley University (online: <http://www.ict4lt.org/>).
- EUROCALL bibliography: <http://www.eurocall-languages.org/resources/bibliography/books.html>
- FITZPATRICK A. and DAVIES G. (eds) (2003), *The Impact of Information and Communications Technologies on the Teaching of Foreign Languages and on the Role of Teachers of Foreign Languages*.
- KARLSSON *et al.* (1995), *Constraint Grammar – A Language-Independent System for Parsing Unrestricted Text*, Mouton de Gruyter.
- WARSCHAUER M. and HEALEY D. (1998), “Computers and language learning: An overview”, in *Language Teaching*, 31: 57-71.
- WARSCHAUER M. (1996), “Computer-assisted language learning: an introduction”, in S. FOTOS (ed), *Multimedia Language Teaching*, Tokyo, Logos International.

# CALL software design principles and the integration of NLP

Jozef Colpaert  
University of Antwerp

## 1. Introduction

The role and shape of solutions for language learning should not be based on a technology-driven (not even NLP-driven) approach, but on an accurate specification of what is needed for a particular language learning situation. We therefore have to create a language learning environment first, defined as an architecture of actors and components and their mutual interactions, before deciding on the language method, media, systems, technologies and NLP routines needed.

Our current research focuses on the implementation of Distributed Language Learning, a conceptual and methodological framework (DLL) for designing language learning solutions in distributed environments.

## 2. Outline of the presentation/demo

In this presentation we will explain the concept of DLL, and show how it has been applied to system development, content structuring, course design and even to the design of a completely new Language Institute.

In the case of system development, we will present a DLL-based software architecture that allows the integration of NLP-routines (Heift and Schulze 2007) on the level of error analysis and answer evaluation but also on the level of the interface (especially for physically, visually or auditory challenged learners). We will show examples in 3D game scripting, language testing and mobile learning.

## 3. Issues and challenges

Our current challenge is the design of an architecture for an intelligent server-based tutoring system for mobile devices, which will be the topic of our next FP7 proposal.

## 4. References

- COLPAERT J. (2007), "Distributed Language Learning", editorial in *Computer Assisted Language Learning*, Vol. 20, No. 1, February 2007: 1-3.
- COLPAERT J. (2007). "Pedagogy-driven design for online language teaching and learning", in *CALICO Journal* 23:3: 477-497.
- COLPAERT J. (2007). "Toward an ontological approach in goal-oriented language courseware design and its implications for technology-independent content structuring", in *Computer Assisted Language Learning*, vol. 19, 2&3: 109-127.
- COLPAERT J. (2004). *Design of Online Interactive Language Courseware: Conceptualization, Specification and Prototyping. Research into the impact of linguistic-didactic functionality on software architecture*. (Doctoral dissertation). University of Antwerp, 2004, 342 p. UMI micropublication number 3141560. Also available on [www.didascaliala.be/doc-design.pdf](http://www.didascaliala.be/doc-design.pdf).
- HEIFT T. and SCHULZE M. (2007), *Errors and Intelligence in Computer-Assisted Language Learning. Parsers and Pedagogues*, Milton Park (Routledge Studies in Computer-Assisted Language Learning (ed. C. Chapelle)).

# CorpusCALL: Challenges and opportunities

Piet Desmet and Hans Paulussen  
K.U.Leuven Campus Kortrijk

## 1. Introduction

This talk is situated in the field of corpusCALL, the use of corpora within CALL (Computer Assisted Language Learning), that has gained growing importance within the CALL research community as can be seen from recent publications (Sinclair 2004, Gavioli 2005, Braun *et al.* 2006, Chambers 2007), and the introduction of SIG communities based on this theme within EuroCall & Calico.

Our research group is quite active within the field of corpusCALL: we have two projects on this domain running at the moment. This should be placed in our general interest in CALL, which has recently led to the foundation of *ALT, Research Center on CALL*. Our current projects involve research topics such as harnessing collective intelligence in e-learning environments, effectiveness of electronic learning platforms, authoring systems for the creation of half-open and open supported tasks and electronic language testing.

## 2. Outline of the presentation/demo

This talk consists of two parts. First of all, we will give an overview of the different approaches of using corpora for foreign language learning, and CALL in particular. The second part will deal with some aspects of corpus creation and the need of standardisation. Both parts will use examples from different projects, including the parallel corpus project REBECA (a collaborative project between K.U.Leuven Campus Kortrijk and FUNDP, Namur) and the recently started DPC project. The Dutch Parallel Corpus project (DPC) is a STEVIN project, organised by a consortium of Dutch and Flemish universities and translation institutes, which aims at compiling a multilingual multifunctional corpus for language technology, translation studies, linguistics and corpusCALL.

Corpora have been created and explored for a long time for different purposes including language technology and linguistics. Only the last ten years has corpus exploration moved to other domains, including foreign language learning and CALL (Computer Assisted Language Learning). Moving away from language specialists to

general language users, corpus exploitation requires an adapted approach which demands different exploitation tools and high quality annotation. We will show where corpora can be useful in language teaching and explain the quality requirements for corpusCALL.

The creation of a corpus has improved considerably over the last ten years, due to the ever growing computer capacity, the interconnectivity between computers of different platform types and the introduction of the internet to the general public. Moreover, most texts are nowadays created electronically. These technological improvements show that corpus creation has become a very easy task, at least as far as collecting text samples is concerned. However, cleaning, structuring and annotating corpora requires careful attention, especially when qualitative exploitation is the ultimate goal. An important improvement in corpus compilation and further exploitation is the introduction of character standardisation (Unicode) and document standardisation via XML (*e.g.* TEI and XCES). Although the XML formats require specific handling, their importance in compiling and distributing text corpora cannot be underestimated. The advantages of XML distribution will be illustrated.

### 3. Issues and challenges

In the context of pure NLP applications (*e.g.* machine translation), corpora are mainly used as linguistic resources to feed a particular application. The corpus is usually transformed into meaningful chunks complying with the requirements of some statistical application. The corpus itself remains invisible to the outside world. In the context of language learning, on the other hand, corpora remain very “visible” to the end-user. Therefore, higher quality standards are required for corpus compilation, annotation and exploitation.

The main challenges in the use of corpora for language learning are situated in further exploitation of corpora. In foreign language learning, corpora can be explored in at least three different stages with reference to the language learning process: (i) corpus extracts in the preparation of language material; (ii) corpus samples during the learning activity and (iii) corpus samples used as feedback after the learning activity. The use of standardised XML formats can improve the exploitation of corpora in each stage.

### 4. References

- BRAUN S., KOHN K. and MUKHERJEE J. (2006), “Corpus technology and language pedagogy”, in *English Corpus Linguistics*, Vol. 3, Frankfurt am Rain, Peter Lang.
- CHAMBERS A. (2007), *Integrating Corpora in Language Learning and Teaching*. Special Issue of *ReCALL*, Volume 17(3).
- DESMET P. and HÉROGUEL A. (2005), “Les enjeux de la création d’un environnement d’apprentissage électronique axé sur la compréhension orale à l’aide du système auteur IDIOMA-TIC”, in *ALSIC (Apprentissage des Langues et Systèmes d’Information et de Communication)*, 8. [http://alsic.u-strasbg.fr/v08/desmet/alsic\\_v08\\_12-poi4.htm](http://alsic.u-strasbg.fr/v08/desmet/alsic_v08_12-poi4.htm)



- DESMET P. (2006), “L’apprentissage/enseignement des langues à l’ère du numérique: tendances récentes et défis”, in *Revue française de linguistique appliquée*, 11: 119-138.
- DESMET P. and EGGERMONT C. (2006), “FRANEL: Un environnement électronique d’apprentissage du français qui intègre des matériaux audio-visuels et qui est à la portée de tous”, in *Cahiers F. Revue de didactique français langue étrangère*, 7: 39-54.
- DESMET P. (2007), “L’apport des TIC à la mise en place d’un dispositif d’apprentissage des langues centré sur l’apprenant”, In *I.T.L. International Journal of Applied Linguistics* 153 (in press).
- DEVILLE G., DUMORTIER L. and PAULUSSEN H. (2004), “Génération de corpus multilingues dans la mise en oeuvre d’un outil en ligne d’aide à la lecture de textes en langue étrangère”, in G. Purnelle, C. Fairon and A. Dister (eds.), *Le poids des mots, Actes des 7es journées internationales d’analyse statistique des données textuelles, JADT 2004*, Louvain-la-Neuve, March 2004: 304-312.
- GAVIOLI L. (2005), *Exploring corpora for ESP learning*, Amsterdam, John Benjamins.
- MACKEN L., TRUSHKINA J., PAULUSSEN H., RURA L., DESMET P. and VANDEWEGHE W. (2007), “Dutch Parallel Corpus: a multilingual annotated corpus”, in *Proceedings of The fourth Corpus Linguistics conference*, University of Birmingham.
- SINCLAIR J. McH. (2004), *How to use corpora in language learning*, Amsterdam, John Benjamins.

Website ALT, Research Center on CALL: [www.kuleuven-kortrijk.be/ALT](http://www.kuleuven-kortrijk.be/ALT)

Website LINGUATIC project: [www.kuleuven-kortrijk.be/linguatic](http://www.kuleuven-kortrijk.be/linguatic)

Website DPC project: [www.kuleuven-kortrijk.be/dpc](http://www.kuleuven-kortrijk.be/dpc)



# The contribution of learner corpus research to TELL

Sylviane Granger  
Université catholique de Louvain

## 1. Introduction

Learner corpus research is a fairly young but highly dynamic research field that emerged in the late 1980s. It focuses on the collection, annotation and computer-aided analysis of vast electronic collections of authentic written and spoken data produced by foreign language learners. The *Centre for English Corpus Linguistics* of the University of Louvain (UCL) has played a key role in shaping the field and demonstrating its tremendous pedagogical potential. In my presentation I will briefly describe the work carried out at Louvain and sketch the numerous possibilities it offers for Technology-Enhanced Language Learning.

## 2. Outline of the presentation/demo

### LEARNER CORPUS COLLECTION

The learner corpora collected at Louvain contain data produced by foreign language learners of English and French. One of the characteristics of the corpora that distinguishes them from other similar collections is that they contain data from a wide range of learner populations. The *International Corpus of Learner English (ICLE)* is a corpus of argumentative essays produced by higher intermediate to advanced learners from 16 different mother tongue backgrounds. The *Louvain International Database of Spoken English Interlanguage (LINDSEI)* is the spoken counterpart of the ICLE and currently covers 11 mother tongue backgrounds. The *French Interlanguage Database (FRIDA)* contains written data from two well-defined learner populations: English- and Dutch-speaking learners of French plus a mixed subcorpus representing a wide range of mother tongue backgrounds. One important characteristic of our corpora is that they are richly documented. In *ICLE* and *LINDSEI* over 20 task and learner variables have been recorded for each of the texts through a detailed profile questionnaire that all learners were requested to complete. All the variables have been stored in a database and can be used by researchers as queries to compile subcorpora that match certain criteria, thus allowing for interesting comparisons (German- vs. Spanish-speaking learners, etc.).

## LEARNER CORPUS ANALYSIS

Two methods of analysis – Contrastive Interlanguage Analysis (CIA) and Computer-Aided Error Analysis (CEA) – have been very popular among learner corpus researchers. CIA is a very powerful, fully automatic heuristic that uncovers the patterns of overuse, underuse and misuse that distinguish learner writing or speech from native or expert user data. CEA is a highly time-consuming but extremely fruitful process that enables researchers to have access to comprehensive catalogues of errors for a given learner population. In my presentation I will illustrate the two methods and describe the systems of error annotation we have designed to make errors amenable to subsequent automated processing.

## PEDAGOGICAL APPLICATIONS

Our learner corpus work has informed a range of pedagogical applications:

- Exercises in the *FreeText* CALL program for learners of French as a Foreign Language (Granger 2003);
- ‘Get it right’ notes and ‘Improve your writing skills’ section in the new edition of the *Macmillan English Dictionary for Advanced Learners* (Rundell and Granger 2007; Gilquin *et al.* in press)
- Web-based error interface: *Exxelant* (Granger *et al.* 2007).

## 3. Issues and challenges

Although learner corpus research has already generated some useful pedagogical applications, its potential is much greater and the extent of learner corpus integration into future applications is likely to increase in the near future. The following are but some of the many avenues for future learner-corpus-informed research: computer-adaptive testing, fleshing out of the descriptors of the Common European Framework levels, error detection and modeling (track what learners have mastered and what they still need to learn; provide appropriate feedback), TELL materials design (notably electronic dictionaries and grammars), incorporation of learner corpus collection and annotation into TELL mobile/web-based environments, speech recognition. All these applications call for synergies between specialists in a wide range of disciplines: TEL, language teaching, language testing, NLP, corpus linguistics and second language acquisition.

## 4. References

- BELZ J.A. and VYATKINA, N. (2005), “Learner Corpus Research and the Development of L2 Pragmatic Competence in Networked Intercultural Language Study: The Case of German Modal Particles”, in *Canadian Modern Language Review* 62.1: 17-48.
- GRANGER S. (ed.) (1998), *Learner English on Computer*, Addison Wesley Longman, London and New York.

- GRANGER S. (2003), "Error-tagged learner corpora and CALL: a promising synergy", in CALICO (special issue on Error Analysis and Error Correction in Computer-Assisted Language Learning) 20(3): 465-480.
- GRANGER S., DAGNEAUX E. and MEUNIER F. (2002), *The International Corpus of Learner English*. CD-ROM and Handbook, Presses universitaires de Louvain, Louvain-la-Neuve. Available from <http://www.i6doc.com>
- GILQUIN G., GRANGER S. and PAQUOT M. (in press), "Learner corpora: the missing link in EAP pedagogy", in *Journal of English for Academic Purposes* 6 (4).
- GRANGER S., HUNG J. & PETCH-TYSON S. (eds.) (2002), *Computer Learner Corpora, Second Language Acquisition and Foreign Language Teaching*, Benjamins, Amsterdam and Philadelphia.
- GRANGER S., KRAIF O., PONTON C., ANTONIADIS G. and ZAMPA V. (2007), "Integrating learner corpora and natural language processing: A crucial step towards reconciling technological sophistication and pedagogical effectiveness", in *ReCALL* 19(3): 252-268.
- RUNDELL M. and GRANGER S. (2007), "From corpora to confidence", in *English Teaching Professional*, 50: 15-18. Reprinted in *MED Magazine*, 46, August 2007 [http://www.macmillandictionary.com/MED-Magazine/August2007/46-Feature\\_CorporatoC.htm](http://www.macmillandictionary.com/MED-Magazine/August2007/46-Feature_CorporatoC.htm).
- NESSELHAUF N. (2004), "Learner corpora and their potential for language teaching", in J. Sinclair (ed.), *How to use corpora in language teaching*, Benjamins, Amsterdam and Philadelphia: 125-152.
- WIBLE D., KUO C.-H., CHIEN F.-Y., LIU A. and TSAO N.-L. (2001). "A web-based EFL writing environment: integrating information for learners, teachers, and researchers", in *Computers and Education*, 37: 297-315.



# Feedback methods in computer assisted pronunciation training applications using automatic speech recognition

Thomas Hansen  
University of Southern Denmark

## 1. Introduction

The field of Computer Assisted Pronunciation Training (CAPT) has seen an explosion in the use of Automatic Speech Recognition (ASR) technology within the past two decades. Contemporary applications come equipped with commercial battle cries of success that leaves one wondering why the use of such applications is not more widespread than it is and also why second language acquisition (SLA) so often still fails.

The main question in this connection is whether the feedback strategies which are presently employed work or how they should be structured to maximize learner benefit.

## 2. Outline of presentation/demo

Contemporary CAPT applications employ a variety of feedback methods, of which the pedagogical value will be discussed. A potential strategy, or roadmap, for improving the effectiveness of ASR in CAPT applications will be outlined for discussion.

## 3. Issues and challenges

Shaping and developing Automatic Speech Recognition technology in such a fashion that a more detailed and constructive level of feedback can be achieved.

## 4. References

HANSEN Th. (2006), "The Four K's of feedback?" In *Proceedings of the 4th International Conference on Multimedia and Information and Communication Technologies in Education (m-ICTE2006)*, Seville: 342-346.

BERNSEN N.O., HANSEN Th., KIILERICH S., MADSEN T. (2006), "Field Evaluation of a Single-Word Pronunciation Training System", in *Proceedings of The Fifth International*

*Conference on Language Resources and Evaluation, LREC 2006*, Genova, May, 2006: 2068-2073.

HINCKS R. (2005), *Computer Support for Learners of Spoken English*, Doctoral dissertation, School of Computer Science and Communication.

HINCKS, R. (2003), "Speech technologies for pronunciation feedback and evaluation", in *ReCall*: 3-20.

NERI A., CUCCHIARINI C. and STRIK H. (2006), "ASR corrective feedback on pronunciation: Does it really work?", in *Proceedings of ICSLP2006*, Pittsburgh: 1982-1985.



# Language technology projects at IDM

Holger Hvelplund  
IDM, Paris

## 1. Introduction

Based on experience with production of electronic versions of monolingual and bilingual dictionaries the emphasis in the presentation will be on how content for cross media products can be produced efficiently, how and why different types of content can/should be integrated; how content can be accessed and adapted in different ways depending on the user, the medium, and the context the content is used in.

## 2. Outline of the presentation/demo

Short demonstration of:

- DPS – a tool for compiling content for different medias and target audiences.
- A few ELT dictionary products with examples of:
  - Different ways of integration language teaching content with the content of the dictionary.
  - How the dictionary can produce content that can be used by language teaching components.
  - Publishing same content on different medias and for different user audiences.
- Teacher resource database.

## 3. Issues and challenges

- Cost efficient production of content for different medias and target audiences.
- Integration of services in new ways where potential in new technologies are exploited. For example, providing language teaching services with Skype, podcast, broadband (like in services from Praxis Language).
- Intelligent user interface with personalized and integrated features.

## 4. References

Recent productions from IDM:

- DPS – Dictionary Production System.
- Production of dictionaries with XDCC for publishers like
  - Oxford University Press (including CD-ROM version of Oxford Advanced Learner's Dictionary).
  - Pearson Education (including CD-ROM and online version of Longman Dictionary of Contemporary English).
  - Macmillan Dictionaries (including CD-ROM version of Macmillan English Dictionary).
  - Cambridge University Press (including CD-ROM and online version of Cambridge Advanced Learner's Dictionary and CD-ROM version of Cambridge Grammar of English).
  - + several others.

# Using corpora in language learning: the Sketch Engine

Adam Kilgarriff  
Lexical Computing Ltd

## 1. Introduction

I am writing an essay about my career plans, and I want to talk about *goals*. How does the word work? What sorts of sentences might I construct around it, with what collocates?

The current range of EFL dictionaries aim to help, and are well-designed, sophisticated tools which specify grammatical patterns and collocates, and show the user a range of example sentences. Often that will be enough. But they are limited to a couple of column inches for a word like *goal* (in which they must cover all of its meanings) and sometimes they just do not cover the case the student is interested in. When that happens, where should they go next?

It is tempting to say that they should go and look in the corpus: after all, that is where the people who wrote the dictionary went. The idea has been discussed at length in the “Teaching and Language Corpora” community. The problem is that reading concordance lines is a skill requiring advanced language competence, and is simply offputting to most learners. The issue may be presented as follows: the dictionary is a highly condensed short summary of the word’s behaviour. The corpus is the raw data for such a summary, not at all condensed or summarized. The user would like a point in between: not as short and minimal as the dictionary, but with a level of abstraction and generalization.

Using techniques from computational linguistics, we have applied just this logic to produce ‘word sketches’ – one-page accounts of the grammatical and collocational behavior of word, as in the figure below.

## 2. Outline of the presentation/demo

We shall demo the Sketch Engine, a tool originally developed for dictionary-making but now being re-engineered for language learners. Recent innovations include a simplified interface, keyword lists, which allow users to find the words which are most

distinctive of a particular subcorpus, and ‘clustered word sketches’, word sketches in which similar words are grouped together.

The word sketch is organized according to grammatical relations, with one list for collocates in each different relation. The relation names (on blue backgrounds) head each list. Collocates are listed according to the grammatical relation they occur in. In contrast to summaries of commonly occurring near neighbours which do not apply grammar, there is no junk: everything is there for an evident linguistic reason.

The first number is the actual number of occurrences of the collocation (taken from the British National Corpus (BNC); all data used here is from the BNC.) The second number is a salience statistic, used for sorting. When working online, the user can click on the number and they are then shown the concordance for the collocation, so if they are unsure what a word is doing in the word sketch, they can promptly find out.

Here, the items are lemmas (dictionary headwords) rather than word forms, so data for *goal* and *goals* are merged. A ‘part of speech tagger’ has been applied to work out, for example, where *post* is a verb (“post the letter”) and where a noun (“goal post”). The word sketch as a whole is for the noun *goal*.

Word sketches were first used for the Macmillan English Dictionary for Advanced Learners, and are now also being used at Oxford University Press, Collins, Chambers Harrap, Le Robert and elsewhere. They changed the way the lexicographers used the corpus. Rather than start with a KWIC concordance for the word, they went straight to the word sketch, as that summarized most of what they needed the concordances for.

Word sketch for *goal* **bnc freq = 10631**

change options

|               |             |            |                  |             |            |                       |            |            |                   |             |            |
|---------------|-------------|------------|------------------|-------------|------------|-----------------------|------------|------------|-------------------|-------------|------------|
| <b>and/or</b> | <b>1112</b> | <b>0.8</b> | <b>object_of</b> | <b>3430</b> | <b>3.1</b> | <b>subject_of</b>     | <b>557</b> | <b>1.0</b> | <b>a_modifier</b> | <b>2546</b> | <b>1.8</b> |
| objective     | 57          | 32.86      | score            | 797         | 75.31      | come                  | 78         | 28.4       | ultimate          | 83          | 42.22      |
| try           | 30          | 32.67      | achieve          | 363         | 48.14      | give                  | 34         | 14.57      | away              | 25          | 32.56      |
| goal          | 32          | 23.39      | concede          | 126         | 47.79      | win                   | 13         | 14.32      | winning           | 31          | 32.56      |
| penalty       | 20          | 22.75      | disallow         | 26          | 34.87      | help                  | 10         | 10.69      | compact           | 34          | 31.79      |
| target        | 22          | 20.1       | pursue           | 75          | 33.13      |                       |            |            | stated            | 17          | 27.88      |
| value         | 33          | 19.36      | attain           | 34          | 29.34      | <b>adj_subject_of</b> | <b>149</b> | <b>1.4</b> | late              | 53          | 27.33      |
| conversion    | 12          | 18.92      | net              | 18          | 26.7       | important             | 10         | 15.32      | dropped           | 11          | 26.98      |
| aim           | 15          | 17.6       | kick             | 36          | 26.2       |                       |            |            | organisational    | 22          | 26.83      |
| mission       | 11          | 16.29      | grab             | 30          | 24.43      |                       |            |            | long-term         | 34          | 25.7       |
| priority      | 10          | 14.13      | reach            | 78          | 23.81      |                       |            |            | common            | 56          | 24.62      |
| strategy      | 11          | 12.28      | set              | 97          | 23.53      |                       |            |            | headed            | 11          | 24.48      |
| point         | 19          | 12.21      | notch            | 10          | 22.81      |                       |            |            | organizational    | 18          | 24.45      |

|                   |             |            |                 |            |            |                   |            |            |
|-------------------|-------------|------------|-----------------|------------|------------|-------------------|------------|------------|
| <b>n_modifier</b> | <b>1181</b> | <b>1.0</b> | <b>modifies</b> | <b>748</b> | <b>0.3</b> | <b>pp_after-p</b> | <b>58</b>  | <b>7.1</b> |
| drop              | 85          | 45.59      | scorer          | 40         | 43.0       | minute            | 37         | 39.18      |
| penalty           | 100         | 45.27      | difference      | 69         | 34.08      |                   |            |            |
| league            | 90          | 37.36      | scoring         | 17         | 29.24      | <b>particle</b>   | <b>86</b>  | <b>4.5</b> |
| consolation       | 24          | 35.39      | Ace             | 18         | 28.33      | back              | 32         | 28.93      |
| opening           | 42          | 31.15      | drought         | 14         | 26.56      | down              | 32         | 28.62      |
| second-half       | 13          | 30.46      | Post            | 34         | 25.55      | up                | 14         | 15.44      |
| first-half        | 12          | 30.04      | Kick            | 17         | 25.19      |                   |            |            |
| minute            | 30          | 21.09      | keeper          | 16         | 24.71      | <b>possessor</b>  | <b>492</b> | <b>4.3</b> |
| half              | 17          | 19.15      | weight          | 21         | 21.01      | England           | 12         | 13.95      |
| policy            | 42          | 18.73      | Lead            | 16         | 20.29      |                   |            |            |
| relationship      | 16          | 13.36      | average         | 10         | 17.56      | <b>pp_from-p</b>  | <b>275</b> | <b>4.1</b> |
| development       | 22          | 13.22      | setting         | 11         | 16.98      | attempt           | 12         | 17.09      |

Goals occur, of course, in sport as well as life. The word sketch highlights the ambiguity. Scanning the ‘object-of’ list, if we *score*, *concede*, *disallow*, *net* or *kick* goals, we are talking sport; if we *achieve*, *pursue*, attain or *reach* them, life. England football fans will be glad to see *England* standing alone in the ‘possessor’ relation to goals!

Word sketches can be explored at <http://www.sketchengine.co.uk> where papers and bibliographical references are also available.

### 3. Issues and challenges

The challenge is to establish the case for corpora in language teaching in general. The second is the attractiveness and usability of the SkE for language learners. The third is of finding or developing appropriate corpora for language learners to use.

### 4. Reference

Kilgarriff A., Rychly P., Smrz P. and Tugwell D. (2004), “The Sketch Engine”, in G. Williams and S. Vessier (eds), *Proceedings of the Eleventh EURALEX International Congress, EURALEX 2004, Lorient, France, July 6-10, 2004*. Lorient: Faculté des Lettres et des Sciences Humaines, Université de Bretagne Sud (Proc EURALEX 2004), Lorient.



# A plurilingual ICALL System for Romance languages

Thomas Koller  
University of Nottingham

## 1. Introduction

The (completed Ph.D.) research described in this abstract deals with the design, development, implementation and evaluation of an interactive plurilingual ICALL (Intelligent Computer-Assisted Language Learning) software system (ESPRIT) for contrastive learning of French, Spanish and Italian. ESPRIT targets learners who are already at an advanced level in at least one of the Romance languages involved. These learners are expected to be familiar with general lexical and grammatical properties of this language. Equivalent properties of the other languages are taught through comparison.

The addressed research questions build upon the general research findings in plurilingual teaching and learning of Romance languages, CALL (Computer-Assisted Language Learning) and ICALL, and the use of animation in language teaching. Formative and summative evaluation processes provided learner assessment data of different components of ESPRIT.

Plurilingual means that grammatical and lexical properties of the languages involved are tightly linked to each other, showing a high degree of similarity in form and function. Plurilingual teaching and learning of Romance languages exploits the similarities between these languages to teach them contrastively and to raise the language awareness of the learner.

The ESPRIT toolset comprises dictionary tools, a concordancer, an input analysis and feedback module, custom-made animated grammar presentations and an authoring tool for animated text. ESPRIT represents a fully functional web-based language learning platform which is designed for autonomous learning. ESPRIT uses a TV metaphor to present language learning materials to the learner. The contents can easily be expanded at any time.

In ESPRIT, learners are free to explore the activities offered and to choose the activities which are of most interest to them. Guided tours, however, provide information and help about which activities form a logical unit, and can be used to suggest in which sequence to work on materials.

## **2. Outline of the presentation/demo**

After providing a short introduction on plurilingual teaching and learning, my presentation will focus on the demonstration of the toolset and the web-based language learning platform developed for ESPRIT. Special emphasis will be put on (a) the NLP techniques integrated into several tools and (b) the applicability of ESPRIT tools and resources to other projects.

Although the similarities between Romance languages have been described extensively in contrastive linguistics for decades, a broader interest in research on plurilingual teaching and learning only emerged in the 1990s. Since then, several European projects have been devoted to plurilingual teaching and learning of Romance languages. The materials developed in these projects do not involve Natural Language Processing (NLP) capabilities and almost exclusively focus on receptive skills. Plurilingual teaching is potentially highly effective, yet plurilingual teaching and learning material is quite hard to obtain. Existing materials only contain a limited amount of reading texts and exercises.

The innovative character of my Ph.D. research lies in the investigation of NLP techniques to enhance the plurilingual teaching and learning of Romance languages. In contrast to existing plurilingual materials, ESPRIT tools allow the learner to work independently on unrestricted learner-retrieved text and to obtain dynamic feedback on learner input. I aimed to develop flexible and interactive, easily expandable software which supports plurilingual teaching and learning of Romance languages and which helps language learners to optimally exploit their existing knowledge in any one Romance language. The single tools and the web-based language learning platform developed for ESPRIT are available at any time on the Internet.

Tools and language data of ESPRIT have already been reused in current projects. Due to their modular character, the tools and language data can easily be integrated in any other project, in which they can be applied to other languages or even language families. Slavic languages, for example, also share a high number of grammatical and lexical properties.

Several ESPRIT tools can also be adapted and provided as Firefox browser plug-ins. As a Firefox extension, an ESPRIT tool would be instantly accessible from any other web page (for example for dictionary look-up). The dictionary tools, lexicon interface components and the concordancer could be adapted as Firefox extensions to provide a wide range of plug-in resources for plurilingual learning of Romance or other languages.

## **3. Issues and challenges**

When learning a third or any further language, learners automatically create links between the properties of currently learned and already learned languages. The plurilingual teaching and learning method, however, is largely unknown to learners



(and teachers). The evaluation for ESPRIT also showed that adult learners varied considerably in the number and type of languages learned already and the degree of fluency therein. Therefore the development of materials for plurilingual teaching and learning posed a number of issues and challenges which differ from second language acquisition and the creation of monolingual language learning materials.

Foreign language teaching in secondary schools and at universities has been largely unaffected by plurilingual research. Language students at both levels only occasionally get the opportunity to learn similar languages simultaneously in a plurilingual setting. As a consequence, it is challenging to identify target learners and to conduct standard institutionalised testing and evaluation of developed plurilingual materials.

The development of plurilingual materials in general has in many cases not been directly connected to research in third language acquisition. Additionally, the majority of existing plurilingual materials tends to be rather descriptive than didactic. Therefore, in my opinion, it would be beneficial for future research in plurilingual teaching and learning to be more tightly linked to research findings of third language acquisition.

#### 4. References

- BLANCHE-BENVENISTE C. (ed.) (1997), *EuRom 4: Metodo de ensino simultâneo das línguas românicas – Metodo para la enseñanza simultánea de las lenguas románicas – Metodo di insegnamento simultaneo delle lingue romanze – Méthode d'enseignement simultané des langues romanes*, Florence, La Nuova Italia Editrice.
- DEGACHE Christian (ed.) (2003), *Intercompréhension en langues romanes. Du développement des compétences de compréhension aux interactions plurilingues, de Galatea à Galanet*, volume 28. Grenoble, LIDILEM, Université Stendhal Grenoble 3.
- KOLLER Th. (2007), *Design, Development, Implementation and Evaluation of a Plurilingual ICALL System for Romance Languages Aimed at Advanced Learners*, PhD thesis, Dublin City University, Dublin.
- MCCANN W.J., KLEIN H.G. and STEGMANN T.D. (2002), *EuroComRom – The Seven Sieves*, volume 5 of Editiones EuroCom. Aachen, Shaker.
- SCHMIDELY J., ALVAR EZQUERRA M. and HERNÁNDEZ GONZÁLEZ C. (2001), *De una a cuatro lenguas. Intercomprensión románica: del español al portugués, al italiano y al francés*. Madrid, Arco Libros.



# Like stars in the firmament: language learning on mobile devices

Agnes Kukulska-Hulme  
The Open University, UK

## 1. Introduction

Mobile learning can be studied as one instance of the ongoing adoption of innovative technologies in the field of education, particularly with a view to understanding learner experience and the potential of the new technologies to transform current practices. Educational uses of mobile technologies offer a rich and complex field of investigation which allows me personally to combine my expertise in e-learning pedagogy with my background in linguistics, language learning, dictionary design and terminology studies (areas I was actively involved in during the 1980s/90s). I'm particularly interested in how mobile devices are changing foreign language learning and how new forms and motivations for language learning might in turn have an effect on attitudes and approaches to multilingual knowledge seeking, global communication and knowledge representation on the web.

My research in mobile learning has been fairly wide-ranging, encompassing studies of how learners read course materials on mobile devices (Waycott and Kukulska-Hulme 2003), surveys of learner-driven mobile innovation (Kukulska-Hulme and Pettit 2006; Pettit and Kukulska-Hulme 2007), critical reviews of evaluation in mobile learning (Traxler and Kukulska-Hulme 2006), reflections on what has been learnt with regard to mobile device usability (Kukulska-Hulme 2007), and issues of collaboration and privacy in contextual learning (Kukulska-Hulme *et al.* 2007). Together with my colleague John Traxler I co-edited the first book on mobile learning to give a coherent account of the field, incorporating a dozen international case studies (Kukulska-Hulme and Traxler 2005). My externally funded projects have also led to the publication of a guide to innovative e-learning with mobile technologies, distributed widely within UK higher and further education. I have tried to make sense of how the field is evolving by studying the possibilities of both formally-designed and user-driven mobile learning (Kukulska-Hulme, Traxler and Pettit 2007). I have also attempted to imagine how mobile language learning will develop (Kukulska-Hulme 2006; Kukulska-Hulme forthcoming).

## 2. Outline of the presentation

Mobile learning is a fast-moving field. It is becoming clearer that device ownership is a factor in adoption, in the type of activity that learners are likely to engage in, and in the integration of learning activities with other aspects of daily life. However, many aspects of mobile learning remain under-explored; for example, connections between mobile and online activity have not yet been investigated in any systematic way, and neither have the implications for the ways that language communication and textual content are used or represented on mobile devices. One interesting observation is that mobile learning projects and initiatives are likely to have outcomes that had not been anticipated by educators or providers of language learning materials, *e.g.* Gilgen (2005) reports that “several new and unexpected uses and results” (p.32) came about in their mobile learning projects.

In my presentation, I will first share and discuss a working classification I have developed of models of participation in mobile language learning: Institutional adoption model; Content delivery model; Proposed activity model; Specified activity model; Content sharing model; and User-generated activity model (Kukulska-Hulme, forthcoming). I’d like to consider the implications for the types of conversations that learners are able to have in these different models of interaction and the extent to which they are able to take part as initiators of learning activities or contributors of language material. If we consider for a moment that each mobile user is interconnected with others, to what extent are they able to be noticed as shining stars in a metaphorical firmament (alluded to in the title of my presentation)? What prospects are there for mobile devices to give learners new opportunities to observe how language is used, request specific types of language support, or share their findings and ideas with others?

## 3. Issues and challenges

*Language captured and shared in context:* Mobile devices are well known for facilitating learning in context – bringing learning closer to real life situations, either spontaneously or in environments designed for those ends. Learners can gather primary data on location, for example by capturing instances of language in use, or use their device to capture and share reflections on language problems and needs, the moment they arise. What are the best ways of fulfilling this potential?

*Fragmented conversations:* Mobile interaction can result in a fragmented experience, especially when use of an online forum is also part of the learning design. We know very little about the learning conversations that mobile devices can, and cannot, facilitate. Ethical and practical issues get in the way of analysing interactions; the specific constraints around mobile learning are still poorly understood.

*Usability:* Often the first thing that people remark on, when they consider mobile learning, is the difficulty of reading from a small screen on a mobile device, or how hard it might be to construct appropriate short text messages within the context of

education. Some then begin to see these difficulties as opportunities or challenges – perhaps shorter texts are better, and have educational value, such as training students to summarize their thoughts or getting teachers to be more precise about instructions. I'm interested in how issues of usability can act in a positive way to instigate reflections on educational goals and changing literacies.

#### 4. References

- CHINNERY G.M. (2006), "Going to the MALL: Mobile Assisted Language Learning", in *Language Learning & Technology*, 10(1): 9-16. Available online: <http://llt.msu.edu/vol10num1/emerging/default.html>
- KUKULSKA-HULME A. and TRAXLER J. (eds) (2005), *Mobile Learning: A Handbook for Educators and Trainers*, Routledge, London.
- KUKULSKA-HULME A. (2006), "Mobile Language Learning Now and in the Future", in Svensson, P. (ed.), *Från vision till praktik: Språkutbildning och Informationsteknik (From vision to practice: language learning and IT)*, Swedish Net University (Nätuniversitetet): 295-310.
- KUKULSKA-HULME A. (2007), "Mobile usability in educational contexts: what have we learnt?", Special issue of the *The International Review of Research in Open and Distance Learning*, 8(2): 1-16. Available online: <http://www.irrodl.org/index.php/irrodl>
- KUKULSKA-HULME A., Traxler J. and Pettit J. (2007), "Designed and User-generated Activity in the Mobile Age", in *Journal of Learning Design*, 2 (1): 52-65. Available online: <http://www.jld.qut.edu.au/>
- PETTIT J. and KUKULSKA-HULME A. (2007,) "Going with the Grain: Mobile Devices in Practice", in *Australasian Journal of Educational Technology (AJET)*, 23 (1): 17-33. Available online: <http://www.ascilite.org.au/ajet/ajet23/ajet23.html>
- SHARPLES M. (2006), *Big Issues in Mobile Learning. Kaleidoscope report*, Available online: <http://mlearning.noe-kaleidoscope.org/repository/BigIssues.pdf>
- WAYCOTT J. and KUKULSKA-HULME A. (2003), "Students' Experiences with PDAs for Reading Course Materials", in *Personal and Ubiquitous Computing*, 7 (1): 30-43.



# Computer-mediated communication for language learning

Marie-Noëlle LAMY  
The Open University, UK

## 1. Introduction

The field of activity captured by the phrase “computer-mediated communication for language learning” recently reached a critical mass, as regards the number of teaching projects and of published papers that have been devoted to it. Chun (2007) suggests that ‘communication’ as used in the phrase ‘computer-assisted communication’ (CMC) receives the most coverage of all topic categories in her overview of recent research based on evidence from two major US journals on technology-mediated language learning, and also comes top of a list of ‘hits’ tracked by one of the two journals in her corpus.

Although caveats are needed due to exclusively US-oriented nature of these results, similar trends are observed in other research cultures, signalling that CMC for language learning (henceforth CMCL) as a field is no longer immature and can be held up to scrutiny. In a volume to be published in November 2007, Lamy and Hampel offer such a scrutiny. The current presentation gives a preview of their findings, and outlines research directions suggested not only by the gaps identified in their study but also by the emergence of new questions raised within neighbouring areas such as multiliteracies research.

## 2. Outline of the presentation

In this section, I present a brief overview of activity in the practice and research of CMCL. I start with the methodological relationships that can be established between CMCL and three related fields, which are: generic (i.e. non-language-oriented) educational CMC, socio-personal CMC, and Computer-Assisted Language Learning or CALL. In the latter field, Warschauer (1995) put CMCL on the map by publishing the first practitioner book on the topic. According to him, the hopes of early adopters of CMCL included giving learners the opportunity to:

- communicate with native speakers and with each other either one-to-one or, more innovatively, one-to-many and many-to-many;

- plan their communication;
- revisit their work, owing to the permanent traces made available to them through the technologies.

I assess whether these expectations been met, and I identify new questions that have arisen along the way, through interrogating the ERIC (Educational Resources Information Centre) database for the years 1992–2005. The results shows that although Warschauer was right to predict a boom in remote communication for language teaching, the hopes of early adopters, as he listed them, were not all fulfilled. Furthermore, a number of meta-studies published since the mid-90s, as well as Warschauer and Kern’s own (2000) review, have shown that the “simple question to which everyone wants an answer ‘Does the use of network-based language teaching lead to better language learning?’ [...] turns out to be not so simple”, and that the CMCL community might do better to abandon the search for improvements in language learning and instead “look to particular *practices of use* including the specifics of learner profiles, task types, process description, discourse, interaction patterns and formal outcomes” (Warschauer and Kern, 2000: 2; original emphasis). By reference to 7 meta-studies, I identify these practices of use and I offer a ‘health check’ of the field, highlighting its achievements, but also the over-coverage of certain topics (e.g. student participation patterns), the under-coverage of others (e.g. assessment) and the recurring concerns expressed by meta-study authors about the quality of research in the CMCL literature.

### 3. Issues and challenges

In this section of the presentation, I identify three types of challenge associated with CMCL. The first one, the unsatisfactory functioning of the practice-research feedback loop, is not specific to this field but can be seen with particularly sharpness in CMCL, perhaps because of the spectacular expansion of the field in a relatively short time.

Major challenge number 2 is at the level of pedagogical practice. I focus on tensions that have not had much exposure in the CMCL research literature so far. Of the two broad groups of tools, text-based and voice-based environments, I argue that both suffer from insufficient attention to the machine-mediated nature of the activity. Many practitioners now agree that online language classes have specific socio-affective or intercultural needs (as CMCL research has indeed been showing) but teachers remain vulnerable to learner disengagement through confusion and overload, because the materiality of the environments is taken as a given rather than being made an explicitly part of the pedagogical considerations that inform teaching design.

Major challenge number 3 is a theoretical and methodological one for researchers. It relates the the multimodal nature of electronic environments and it concerns the choice of theoretical and methodological frameworks for the analysis of learner conversations, when such conversations may be carried out via a range of interrelating semiotic systems, some wired into the machine and others freely deployed or created by the learners, many of them of a non-linguistic nature.



Finally, I show how the nature of challenges 2 and 3 point to the importance of multiliteracies research for the future of CMCL.

#### 4. References

- CHUN D.M. (2007), "Come Ride the Wave: But Where is it Taking Us?", in *The CALICO Journal* 24(2): 239-252.
- HASSAN X., HAUGER D., NYE G. and SMITH P. (2005), "The Use and Effectiveness of Synchronous Audiographic Conferencing in Modern Language Teaching and Learning (Online Language Tuition): A Systematic Review of Available Research", in *Research Evidence in Education Library*, London: EPPI-Centre, Social Science Research Unit, Institute of Education, University of London.
- HUBBARD P. (2005), "A Review of Characteristics in CALL Research", in *Computer Assisted Language Learning* 18(5): 351-368.
- JUNG U. (2005), "CALL – Past, Present and Future: A Bibliometric Approach", in *ReCALL* 17(1): 4-17.
- KERN R. G. (2006), « La Communication médiatisée par ordinateur en langues: recherches et applications récentes aux USA », in F. Mangenot and C. Dejean-Thircuir (eds), *Les Echanges en ligne dans l'apprentissage et la formation, le français dans le monde, recherches et applications* 40: 17-29.
- LEVY M. (2000), "Scope, Goals and Methods in CALL Research: Questions of Coherence and Autonomy", in *ReCALL* 12(2): 170-195.
- LIU M., MOORE Z., GRAHAM and LEE, S. (2002), "A Look at the Research on Computer-based Technology Use in Second-language Learning: A Review of the Literature from 1990-2000", in *Journal of Research on Technology in Education* 34(3): 250-273.
- LAMY M.-N. and HAMPEL R. (2007 forthcoming)? *Online Communication for Language Teaching and Learning*, Palgrave Basingstoke, Palgrave MacMillan.
- WARSCHAUER M. (1995), *Virtual Connections: Online Activities and Projects for Networking Language Learners*, Honolulu: Second Language Teaching and Curriculum Centre, University of Hawaii.
- WARSCHAUER M. and KERN R. (eds) (2000), *Network-based Language Teaching: Concepts and Practice*, Cambridge, Cambridge University Press.
- ZHAO Y. (2003), "Recent Developments in Technology and Language Learning: A Literature Review and Meta-analysis", in *The CALICO Journal* 21(1): 7-27.



# Writing English as a second language: A proofreading tool

Claudia Leacock  
The Butler Hill Group

## 1. Introduction

This work combines Natural Language Processing (NLP) and Machine Learning (ML) techniques for detecting and correcting grammatical errors in the writing of English Language Learners (ELL).

The *Writing English as a Second Language* tool being developed at Microsoft Research (for which the presenter is a consultant) focuses on those areas of grammar that pose special challenges for English language learners. This presentation focuses on those problems that are hardest for ELLs – the use of determiners and of prepositions – although the system also identifies gerund/infinitive confusion, auxiliary verb presence and choice, over-regularized verb inflection (*writed vs. wrote*), adjective/noun confusion (*China book vs. Chinese book*), word order errors, and mass vs. count noun errors (*much knowledge vs. many knowledges*).

## 2. Outline of presentation

I will describe the system's major components:

1. The *Suggestion Provider* consists of Individual *error identification modules* that identify potential errors. These modules are flexible and can identify errors using ML techniques rules, regular expressions or a combination of the three.

For the preposition and determiner correction modules, a classifier is used that is trained on edited native English. For each potential insertion point of a determiner or preposition in that training data, a vector of features is extracted from the context.

2. The *Language Model*, which selects the most likely suggestion(s), is a 5-gram model trained on the English gigaword corpus.
3. The *Example Provider* retrieves relevant example sentences from the web to help the user select the most appropriate rewrite. This innovative component generates

an exact string query including a window of context around the suggested correction. The query is issued to a search engine, and the retrieved sentences are ranked and presented to the user.

For system accuracy, we will present two different evaluations: (1) Automatic evaluation on copyedited native text – under the assumption that it contains no errors. (2) Human evaluation of essays written by Chinese ELLs.

### 3. Issues and challenges

Given the very high frequency with which prepositions and determiners occur in English, the false flag rate must be very low in order to be acceptable. While a native speaker can easily identify and ignore a false flag, language learners would have to take time inspecting the example sentences to decide which is correct – and even then may get confused. The system currently uses handcrafted heuristics to minimize false flags – which is laborious and requires retuning each time a model is retrained. The system’s developers are investigating a learned ranker to replace the handcrafted heuristics.

Traditionally, grammar checkers give the writer a mini-multiple-choice test when providing potential corrections. This is fine for native writers – but simply poses another challenge for language learners. The Example Provider is an innovative method for enabling the user to make an informed decision. Initial response has been positive but will require refinements based on the results of extensive user testing.

### 4. References

- BURSTEIN J. and LEACOCK C. (eds) (2006), *Natural Language Engineering: Special Issue on Using NLP in Educational Applications*, 12:2.
- BURSTEIN J., CHODOROW M. and LEACOCK C. (2004), “Automated Essay Evaluation: The Criterion Online Writing Service”, in *AI Magazine* 25:3: 27-36.
- BURSTEIN J., CHODOROW M. and LEACOCK C. (2003), “Criterion Online Essay Evaluation: An Application for Automated Evaluation of Student Essays”, in *Proceedings of the Fifteenth Annual Conference on Innovative Applications of Artificial Intelligence Conference (IAAI-03)*, Acapulco.
- CHODOROW M., TETREault J.R. and HAN N.-R. (2007), “Detection of Grammatical Errors Involving Prepositions”, in *Proceedings of the 4<sup>th</sup> ACL-SIGSEM Workshop on Prepositions*: 25-30.
- CHODOROW M. and LEACOCK C. (2000), “An Unsupervised Method for Detecting Grammatical Errors”, in *Proceedings of the 1<sup>st</sup> Annual Meeting of the North American Chapter of the Association for Computational Linguistics*, Seattle, WA.
- HAN N.-R., CHODOROW M. and LEACOCK C. (2006). “Detecting Errors in English Article Usage by Non-Native Speakers », in *Natural Language Engineering* 12:2.
- HAN N.-R., CHODOROW M. and LEACOCK C. (2004). “Detecting Errors in English Article Usage with a Maximum Entropy Classifier Trained on a Large, Diverse Corpus”, in *Proceedings of the 4th International Conference on Language Resources and Evaluation*. Lisbon.

- IZUMI E., UCHIMOTO K., SAIGA T., SUPNITHI T. and ISAHARA H. (2003), "Automatic Error Detection in the Japanese Learners' English spoken Data", in *The Companion Volume to the Proceedings of 41st Annual Meeting of the Association for Computational Linguistics, July 2003*: 145-148.
- LEACOCK C., and CHODOROW M. (2003), "Automated Grammatical Error Detection", in M.D. Shermis and J. Burstein (eds), *Automated Essay Scoring: A Cross-Disciplinary Perspective*. Hillsdale, NJ, Lawrence Erlbaum : 195-208.
- LEACOCK C., and CHODOROW M. (2001), "A Corpus-Based Approach to Diagnosing Grammatical Errors", in *Corpus Linguistics and Language Teaching Conference*, Boston, MA.
- LIU T., ZHOU M., GAO J., XUN E. and HYAN C. (2000), "PENS: A Machine-Aided English Writing System for Chinese users", in *Proceedings of ACL 2000*: 529-536.
- SHERMIS M.D., BURSTEIN J. and LEACOCK C. (2006), "Applications of Computers in Assessment and Analysis of Writing", in C.A. McArthur, S. Graham and J. Fitzgerald (eds), *Handbook of Writing Research*, New York, Guilford Press.
- TURNER J. and CHARNIAK E. (2007), "Language Modeling for Determiner Selection", in *Human Language Technologies 2007: The Conference of the NAACL; Companion Volume, Short Papers*: 177-180.



# Second language acquisition theory and TELL

Fanny Meunier  
Université Catholique de Louvain

## 1. Introduction

As a researcher in Second Language Acquisition (and more specifically instructed second language acquisition) and a teacher of English as a foreign language, my aim in this presentation is twofold: first, stress the importance of the two 'L's in technology-enhanced language learning; and secondly, address the convergences and divergences that exist between big issues in second language acquisition and TELL.

## 2. Outline of the presentation

The appeal of new technologies in language learning has undoubtedly played a role in downgrading the focus on teaching/learning methodologies per se, be they technology-enhanced or not. Wible (2005:2) even states that whilst massive resources are invested in the development of information technology for e-learning, little concentrated effort is devoted to bridging the gap between technology and second/foreign language education. I will argue that an additional gap still increases the complexity of the situation, i.e. the gap that exists between SLA theory and second/foreign language education.

According to Wible again (2005:72), TELL would benefit from a shift from what the technology is capable of doing to what the learner actually needs. I will demonstrate in my presentation that such a shift ideally requires insights into SLA theory.

The issues that will be dealt with include the positive or negative influence that technology may have in addressing some internal/external variables in SLA (L1 background, learners' characteristics, learning styles, input, etc.), the receptive/productive skills dichotomy, a number of cognitive aspects (awareness, saliency, elaboration, rehearsal, etc.) and types of feedback.

## 3. Issues and challenges

First, closer collaboration should be encouraged between the technological, acquisitional and methodological paradigms of language learning. Secondly, TELL should consider what I call (see Meunier forthcoming) 'principled eclecticism in

learning' as one of its future challenges. Learners in TELL environments should have access to observational, descriptive or explanatory options, together with opportunities for immediate feedback. Third, learnability issues (defined here as the input/output efficiency of some method or approach) should become more central in order to validate the efficiency of TELL.

#### 4. References

- CHAPELLE C. (1998), "Multimedia CALL: Lessons to be Learned from Research on Instructed SLA", in *Language Learning and Technology*, 2(1): 22-34.
- CHAPELLE C. (2003), *English Language Learning and Technology*, Benjamins, Amsterdam and Philadelphia.
- KASPER L. (2000), "New technologies, new literacies: focus discipline research and ESL learning communities", in *Language Learning & Technology*, vol. 4, n°. 2: 105-128.
- LIGHTBOWN P. and SPADA N. (2003) *Factors Affecting Second Language Learning. How Languages Are Learned*, Revised edition, Oxford University Press, Oxford.
- MEUNIER, F. (forthcoming) "Corpora, cognition and pedagogical grammars: An account of convergences and divergences", in S. De Knop and T. De Rycker (eds), *Cognitive Approaches to Pedagogical Grammar*, Mouton de Gruyter, Berlin.
- WIBLE D. (2005), *Language Learning and Language Technology: Toward Foundations for Interdisciplinary Collaboration*, Crane, Taipei.
- WIBLE D. (in press), "Multiword Expressions and the Digital Turn", in F. Meunier and S. Granger (eds), *Phraseology in Foreign Language Learning and Teaching*, John Benjamins, Amsterdam.



# Detecting syntactic interference

John Nerbonne  
University of Groningen

## Introduction

This presentation involves joint work with Wybo Wiersma (Groningen) and Timo Lauttamus (Oulu). It applies techniques from quantitative computational linguistics to the problem of detecting frequent effects of first language interference in second language learning. We focus on production interference in syntax.

Second language learners typically differ syntactically from native speakers not only in making outright errors, but also in overusing and under-using some constructions – all of which we subsume under INTERFERENCE. In approaching the phenomenon of interference computationally, we were motivated both to attempt to identify interference effects more systematically, and also to attempt to quantify a level of aggregate interference, a goal Weinreich (1953: 63) found worthwhile, but which he speculated to be unreachable:

No easy way of measuring or characterizing the total impact of one language on another in the speech of bilinguals has been, or probably can be devised. The only possible procedure is to describe the various forms of interference and to tabulate their frequency.

## 2. Outline

Following a suggestion by Aarts and Granger (1998), we model the syntax of the second-language learners of English via the parts of speech (POS) they use and the sequences in which the POS appear. The idea is to compare the POS sequences used by second-language learners to those used by natives. Concretely, we examine the distribution of triplets of POS in a large corpus of English as used by adult Finnish immigrants to Australia, and compare this distribution to that of their children, who immigrated as children and speak English at a near-native level. To assay the “total impact” as Weinreich wished, we examine the differences between the two distributions via a permutation test, which is implemented in a Monte Carlo fashion. To identify systematically the areas of difference, we examine the frequent POS triples that contribute most to the overall differences in the two distributions.

### 3. Issues and challenges

Although we can assay “total impact” and also identify areas of syntactic differences straightforwardly, there are several points at which improvement would be useful and interesting. First, the technique works at a high level of aggregation. This is not a serious problem in language contact study, which is interested in exactly such population effects, but it is a problem if one wishes to analyse the work of individual second-language learners. Second, and related, we should wish to study the influence of individual speakers on the approach, as Sanders (2007) has. Third, our approach assumes that POS sequences represent syntax well (for the purpose of assaying differences), an assumption which is justified by the endocentricity of syntax, but this assumption could also be tested. Finally, we need to analyse more data sets involving different languages since the Finnish effects on English may be a special case.

We are also interested in receptive interference (Moberg et al. 2007), but will not have the time to present that work here.

### 4. References

- AARTS J. and GRANGER S. (1998), “Tag sequences in learner corpora: A Key to Interlanguage Grammar and Discourse”, in S. Granger (ed), *Learner English on Computer*, London, Longman: 132-141.
- MOBERG J., GOOSKENS C., NERBONNE J. and VAILLETTE N. (ca. 2007), “Conditional Entropy Measures Intelligibility among Related Languages”, accepted to appear in F. Van Eynde, P. Dirix, I. Schuurman and V. Vandeghinste (eds.) *Proceedings of Computational Linguistics in the Netherlands 2006*, Amsterdam, Rodopi.
- NERBONNE J. and WIERSMA W. (2006), “A Measure of Aggregate Syntactic Distance”, in J. Nerbonne and E. Hinrichs (eds), *Linguistic Distances*, Workshop at the joint conference of International Committee on Computational Linguistics and the Association for Computational Linguistics, Sydney, July, 2006: 82-90.
- SANDERS N. (2007) “Measuring Syntactic Difference in British English”, in *Proceedings of the ACL 2007 Student Research Workshop*, Prague: ACL: 1-6.
- WEINREICH, Uriel ([1953], 1968, 1974), *Languages in Contact*. The Hague, Mouton.

# Planning a smart phone system to support self-directed L2 vocabulary learning

Richard Pemberton  
University of Nottingham

## 1. Introduction

There are two major problems with vocabulary learning that almost every language learner will be familiar with:

- learning enough frequent vocabulary to be able to read and listen fluently;
- retaining the vocabulary that we have learned.

The first problem involves a considerable amount of time. To take English for example, in order to be able to understand unsimplified texts, you need to know some 3,000–4,000 of the most common English word families (Nation & Waring 1997; Nation 2001). The figure is likely to be upwards of 5,000 word families if fluent reading for pleasure is the aim.

Equally, if vocabulary is to be retained, the learner needs to spend a lot of time in conscious processing or repeating of the target items (explicit learning) and/or in extensive language use (implicit learning). These problems of time are of course even worse for the busy adult learner living outside the target country.

One type of technology which has the potential to save time on the go and to support both implicit and explicit learning is the mobile phone. However, recent mobile phone systems supporting vocabulary learning have tended to use one medium only – e.g. text messages (Pincas 2004; Song & Fox 2005), e-mail (Thornton & Houser 2005) or images (Joseph et al 2005) – and to have involved designed rather than learner-located, learner-generated and learner-shared materials and activities.

## 2. Outline of presentation

In this presentation I will first propose a smart phone system that could use the phone's full capabilities (see e.g. Kukulska-Hulme & Shield 2007: 20) to support self-directed vocabulary learning.

I will then exemplify and discuss the desirability of various potential features (both 'existing' and 'to-be-created'), including:

- the creation of subtitles for video recordings (cf. Fallahkhair et al 2007);
- OCR scanning of camera shots of written text;
- the use of a personalisable wordlist/dictionary with testing functions;
- a testing ‘scheduler’ based on the principles of spaced repetition and expanding rehearsal (Ellis 1995);
- the use of a ‘producing’ dictionary with associative functions;
- geotagging;
- a simple advising system to provide guidance re vocabulary learning.

### 3. Issues and challenges

There are a number of issues to be clarified and challenges to be met, including:

- How long will it take before more phone models provide enough screen size for subtitled videos?
- Would it be too expensive for Europe? Which countries would it be best to target?

### 4. References

- ELLIS N.C. (1995), “The psychology of foreign language vocabulary acquisition: implications for CALL”, in *Computer Assisted Language Learning* 8(2-3): 103-128.
- FALLAHKHAIR S., PEMBERTON L. and GRIFFITHS R. (2007), “Development of a cross-platform ubiquitous language learning service via mobile phone and interactive television”, in *Journal of Computer Assisted Learning* 23: 312-325.
- JOSEPH S., BINSTED K. and SUTHERS D. (2005), “PhotoStudy: vocabulary learning and collaboration on fixed and mobile devices”, in H. Ogata, M. Sharples, G. Kinshuk and Y. Yano (eds), *Proceedings of the third IEEE International Workshop on Wireless and Mobile Technologies in Education 2005*, Los Alamitos, CA: IEEE: 206-210.
- KUKULSKA-HULME, A. and SHIELD L. (2007), “An overview of mobile assisted language learning: can mobile devices support collaborative practice in speaking and listening?”, in EuroCALL 2007. <http://vsportal2007.googlepages.com/collaborativepractice> [Accessed 15 September 2007]
- KUKULSKA-HULME A., TRAXLER J. and PETTIT J. (2007), “Designed and user-generated activity in the mobile age”, in *Journal of Learning Design* 2(1): 52-65. <http://www.jld.qut.edu.au> [Accessed 15 September 2007.]
- NATION I.S.P. (2001), *Learning Vocabulary in Another Language*, Cambridge, Cambridge University Press.
- NATION P. and WARING R. (1997), “Vocabulary size, text coverage and word lists », in N. Schmitt and M. McCarthy (eds), *Vocabulary: description, acquisition and pedagogy*, Cambridge, Cambridge University Press: 6-19.
- PINCAS A. (2004), “Using mobile phone support for use of Greek during the Olympic Games 2004”, in *International Journal of Instructional Technology & Distance Learning* [http://www.itdl.org/Journal/Jun\\_04/article01.htm](http://www.itdl.org/Journal/Jun_04/article01.htm) [Accessed 19 September 2007.]
- SONG Y. and FOX R. (2005), “Integrating web-based ESL vocabulary learning for working adult learners”, in H.Ogata, M. Sharples, G. Kinshuk and Y. Yano (eds), *Proceedings of the third IEEE International Workshop on Wireless and Mobile Technologies in Education 2005*, Los Alamitos, CA, IEEE: 154-158.
- THORNTON P. and HOUSER C. (2005), “Using mobile phones in English education in Japan”, in *Journal of Computer Assisted Learning* 21: 217-228.

THORNTON P. and SHARPLES M. (2005), "Patterns of technology use in self-directed Japanese language learning projects and implications for new mobile support tools". In H.Ogata, M. Sharples, G. Kinshuk and Y. Yano (eds), *Proceedings of the third IEEE International Workshop on Wireless and Mobile Technologies in Education 2005*, Los Alamitos, CA, IEEE: 203-205.



# The dictionary of the future

Michael Rundell

Lexicography MasterClass Ltd and Macmillan Dictionaries

## 1. Introduction

Donald Rumsfeld's famous reflections on "what we know we don't know and what we don't know we don't know" apply to most forms of futurology, and certainly to any attempt to predict what might happen in the world of reference materials. This talk is at the interface of language-learning, lexicography, NLP, and delivery media, and will outline some possible future directions for dictionaries aimed at learners of English.

## 2. Outline of the presentation/demo

The monolingual learner's dictionary (MLD) follows a model whose essential characteristics were developed in the 1930s (by people like Harold Palmer and Michael West) and found concrete form in A.S Hornby's ur-MLD, published in 1942. This was later to morph into the *Oxford Advanced Learner's Dictionary*, which is still going strong after 60 years. There have been two big changes in the intervening period:

- the arrival of corpora in the 1980s, which led to great improvements in quality as lexicographers got access to objective language data (Sinclair 1987);
- the growth of competition: Hornby's dictionary had the field to itself till 1978, but four other contenders have since entered the fray, and this has helped to drive innovation (Rundell 1998).

Though content and accessibility of these dictionaries has steadily improved, the basic model hasn't fundamentally altered. But like any other kind of reference resource, the MLD can't fail to be affected by the biggest change of all – the arrival of the Web.

I will look first at signs that the old model is beginning to break down (for example, the fact that electronic versions of MLDs have begun to include content not present in the print editions); then consider current challenges and opportunities; and finally suggest what the MLD might look like ten years from now.

### 3. Issues and challenges

Probably the most obvious point is that the MLD may no longer be quite ‘M’ or ‘D’. The old binary choices in reference publishing (monolingual or bilingual, dictionary or encyclopedia, advanced or intermediate) may no longer be relevant. Customization and personalization are likely new directions, so the current globally-marketed one-size-fits-all package will probably be unpicked. From the point of view of *content*, lexicographers and linguists have never been better placed. The age of data-sparseness is behind us, and we have fantastic language resources at our disposal (corpora of infinite size, and language-analysis software of increasing power and sophistication: e.g. Kilgarriff and Rundell 2002). Effective exploitation (in dictionaries) of learner corpora has only just begun (Rundell and Granger 2007; Gilquin, Granger and Paquot forthcoming), and there is much more to be done on this front. Essentially, we can do anything, and there are plenty of areas of the language that dictionaries do not yet describe adequately.

The challenges include:

- matching content to users’ needs: one of the issues here is that so much reference material is available at no cost on the Web (Google, Wikipedia etc.), so we have to be clear about what to focus on. The challenge is to work out how to provide information which learners need, and which is *either* not available elsewhere, *or* not available in an easy-to-use form that takes account of learners’ needs (and limitations);
- how reference data will be delivered: the current platform for electronic dictionaries is the CD-ROM, already an ageing technology with obvious limitations, so what (in addition to online access) might replace it? A possible scenario is to see our reference materials as a set of components which customers can mix and match according to their needs. For example, a learner from China doing a Masters in agriculture at a British university could have access to a ‘core’ ELT dictionary with the option of Chinese translations, academic-writing aids, and subject-specific terminology. The resources thus become less static, more dynamic;
- the hardest question: how to fund all this development? Electronic versions of MLDs have been around for 15 years or so, but none have yet made any money (and they cost a lot to develop). New revenue models need to emerge, and these could include advertising. (The typical users of MLDs are an adman’s dream: young, intelligent, aspirational etc.). To be discussed...

### 4. References

- DE SCHRYVER G.-M. (2003), “Lexicographers’ dreams in the electronic dictionary age”, in *International Journal of Lexicography*, 16.2: 143-199.
- GILQUIN G., GRANGER S. and PAQUOT M. (forthcoming), “Learner corpora: the missing link in EAP pedagogy », in *Journal of English for Academic Purposes*.



- KILGARRIFF A. (2006), "Collocationality and how to measure it", in E. Corino, C. Marelo, C. Onesti (eds.), *Proceedings of 12th EURALEX International Congress*, Alessandria, Edizioni Dell'Orso: 997-1004.
- KILGARRIFF A. and RUNDELL M. (2002), "Lexical Profiling Software and its Lexicographic Applications – a Case Study", in A. Braasch and C. Povlsen (eds.), *Proceedings of the 10th EURALEX Copenhagen Proceedings*: 807-818.
- RUNDELL M. (1998), « Recent trends in English pedagogical lexicography », in *International Journal of Lexicography*, 11.4: 315-342.
- RUNDELL M. and GRANGER S. (2007), "From corpora to confidence", in *English Teaching professional*, 50 : 15-18.
- SINCLAIR J.M. (ed.) (1987), *Looking Up: the COBUILD project in lexical computing*, London, Collins ELT.



# Macmillan English Campus: A case study

Emma Shercliff  
MacMillan English Campus

## 1. Introduction

Macmillan English Campus is an online practice environment designed for the learning and teaching of English as a Foreign Language. It was developed in conjunction with one of the world's leading language schools, Cultura Inglesa, Sao Paulo, and is today being used by over 90,000 students worldwide.

Macmillan English Campus consists of two components:

1. A flexible database of over 3,000 highly interactive language activities, developed by Macmillan's leading ELT authors. These activities include interactive language exercises, listening tasks, pronunciation exercises, vocabulary exercises, progress tests, exam preparation exercises, language games, web projects and weekly news items. All users also have access to an online version of the Macmillan English Dictionary.
2. Sophisticated content management software, allowing institutions to manage their users and chart our online resources to their own courses and course materials. The Macmillan English Campus platform includes an electronic mark book and personalisation tools for each user.

The concept behind the Macmillan English Campus is that language learning can be greatly enhanced by an effective combination of face-to-face teaching and customized online support materials. It is this blended learning solution that makes the Macmillan English Campus unique. It ensures that our users continue to receive face-to-face tuition and contact with their teachers whilst remaining free to study online within a controlled learning environment.

## 2. Outline of the presentation/demo

Macmillan English Campus, an online language learning environment, is at the cutting edge of TEL developments. It makes use of exciting new web technologies to deliver an innovative learning experience for both teachers and students. As such, it provides an excellent example of the challenges encountered by practitioners on a daily basis when exploring new and more effective ways of dealing with language.

This presentation will comprise a case study of the Macmillan English Campus, which was launched in 2003. I will outline the concept behind the ‘blended learning’ pedagogy of the Macmillan English Campus and will then address the issues and challenges we have faced over the past four years, with specific reference to the experience of language teachers wishing to integrate technology-enhanced learning into their teaching programmes for the first time. I will outline the enhancements we have incorporated into the Macmillan English Campus learning platform as a result of user feedback and outline future developments we have planned for 2008 and beyond. In the light of our extensive experience developing online learning applications, I will also highlight what we believe to be the limitations of technology in a language learning context.

I will give specific examples of TEL methods and tools and demonstrate some of the new functionality, such as teacher-to-student messaging, recently incorporated into the Macmillan English Campus.

The Macmillan English Campus has been adopted by a number of teaching prestigious institutions, schools and universities worldwide, including the International House World Organisation, the British Council and the Bell Schools network. The presentation will draw on Macmillan English Campus’s widespread experience in the field and is intended to focus on practice rather than theory. Much of our publishing is driven by user responses to our learning platform and we have therefore developed sophisticated mechanisms for gathering and evaluating feedback from users across five continents.

By sharing the experiences of Macmillan English Campus, I will offer a practical insight into the challenges of developing materials for technology enhanced language learning which I hope will stimulate comment and debate.

### **3. Issues and challenges**

The issues and challenges faced at Macmillan English Campus include the following:

#### **Publishing**

- How to publish for the web – inventing a new authoring tool
- Ensuring the English Campus is tailored to a school’s teaching programmes and pedagogic style to enable true ‘blended learning’
- Digital asset management
- The limitations of technology *e.g.* speech recognition tools

#### **Commercial**

- Investment: enormous cost of developing online platform
- Perception amongst certain customers that digital product should be free
- Initial assumption amongst certain customers that teachers could author their own material at lower cost

- Decision making process slow: initial reluctance of institutions to embrace online learning platform as adoption necessarily involves a change in pedagogy and teaching methodology

#### Training & Support

- Teacher training
- Lack of specific technical expertise within language learning organisations
- Creating an online community for Macmillan English Campus users to share examples of best practice (daily blog now available at [www.macmillanenglishcampus.com/support](http://www.macmillanenglishcampus.com/support))



# **The *Base lexicale du français (BLF)*, a free web-based learning environment for French vocabulary**

Serge Verlinde  
Katholieke Universiteit Leuven

## **1. Introduction**

This presentation deals with recent developments in web-based electronic learner's dictionaries and their use in CALL (computer assisted language learning) applications.

My research interests are the lexicon and its structure, corpus analysis and CALL. I am coauthor of the *Dictionnaire d'apprentissage du français des affaires* (DAFA) and have developed, together with Thierry Selva, the *Base lexicale du français* ([www.kuleuven.be/ilt/blf](http://www.kuleuven.be/ilt/blf)), a free accessible learning environment (online dictionary and exercises) for French vocabulary.

## **2. Outline of the presentation/demo**

Today electronic dictionaries, and electronic pedagogical dictionaries in particular, are much more than an electronic version of a paper dictionary. Combined with NLP applications, they may be turned into a powerful language teaching/learning (and research) tool.

The BLF is a free web-based learning environment of a new generation, which combines

- a learner's dictionary or lexical database (*Dictionnaire d'apprentissage du français langue étrangère ou seconde – DAFLES*);
- a corpus of newspaper texts;
- a CALL application (Alfalex);
- direct access to other freely accessible lexical resources on the web.

The BLF was developed from scratch and it is based on a relational database.

Theoretically, and as far as the data-structuring allows, one should be able to launch any query on a lexical database. These queries could apply to the lexicon itself (nomenclature, word combinations) as well as to its properties (*e.g.* grammatical category for the nomenclature; lexical function for the word combinations), and even

to a series of characters contained in a cell of the database (*e.g.* a query concerning all definitions encompassing the noun *action* or all verbs used with a prepositional group introduced by the preposition *à*). In the BLF, we have tried to reach these goals within the didactic perspective of teaching/learning French as a foreign language, as well as exploit these resources for research purposes.

The corpus is used to provide both examples of the use of multiword units (word combinations) and sentences for the exercises in the CALL application (ALFALEX) by using NLP-tools.

ALFALEX offers about ten different types of exercises relating to 'the words' most important features:

- formal features (morphology, verb conjugation, derivation);
- intrinsic features (gender);
- combinatorial features (use of prepositions after verbs, nouns and adjectives, multiword units);
- lexical relations (synonyms, *schémas actanciels* or words encountered in the same communicative situation: *e.g.* how do we designate the act of killing (*assassiner*) a person ? *un assassinat* ; what do we call the person who killed another person ? *un assassin* ; and the person who was killed? *la victime*);
- translation (decoding: French > Dutch, encoding: Dutch > French).

The exercises listed above are semi-automatically generated through direct use of the information in the lexical database and the corpus (for the contextual exercises). Directional and constructive feedback is provided: by means of hyperlinks, the user can access the lexicographical description, which is available for almost every item in the exercises. Twice a year, ALFALEX also automatically generates a qualitative report for every user of the environment.

Unfortunately, non-commercial dictionaries such as the DAFLES only cover a part of the lexicon. Therefore, if a user submits a word which is not listed, he will be redirected to other lexical resources available on the internet. He also has access to other free web resources (*e.g.* corpora, semantic networks) for French.

### 3. Issues and challenges

- How do learners use electronic dictionaries and the most recent resources available on the web? How can the dictionary be tailored to users' real needs? Tracking and logging the real use of (web-based) electronic dictionaries is certainly the first step of discovering this.
- Learners still have problems using the dictionary for decoding purposes (length and structure of the entries, incomplete lexical description, identification of multiword



units, ...). Could NLP applications (e. g. the use of a parser), combined with a dictionary/lexical database, be helpful?

-- Encoding is even more complicated. How can we develop a real writing assistant?

#### 4. References

- ABEL A. and WEBER V. (2000), ELDIT – A prototype of an innovative dictionary, in U. Heid *et al.* (eds), *Proceedings EURALEX, The Ninth EURALEX International Congress*, Stuttgart : 807-818.
- ALDABE I., ARRIETA B., DÍAZ DE ILARAZZA A., MARITXALAR M., NIEBLA I., ORONOZ M. and URIA L. (2006), “The use of NLP tools for Basque in a multiple user CALL environment and its feedback”, in P. Mertens, C. Fairon, A. Dister and P. Watrin (eds), *Verbum ex machina. Actes de la 13e conférence sur le Traitement automatique des langues naturelles*, Louvain-la-Neuve : 815-824 (Cahiers du Cental 2).
- ANTONIADIS G., ECHINARD S., KRAIF O., LEBARBÉ T. and PONTON C. (2005), « Modélisation de l'intégration de ressources TAL pour l'apprentissage des langues : la plateforme MIRTO », in *Alsic.org*, vol. 8: 65-79.
- ANTONIADIS G., ECHINARD S., KRAIF O., LEBARBÉ T., LOISEAU M. and PONTON C. (2004), “CALL: from current problems to NLP solutions MIRTO: a user-oriented NLP based teaching platform”, in *Proceedings of EuroCALL Conference 2004*, Vienna.
- BLUMENTHAL P. (2006): *Wortprofil im Französischen*, Tübingen.
- BLUMENTHAL P. and HAUSMANN F.J. (eds.) (2006), *Collocations, corpus, dictionnaires*, in *Langue française*, 150.
- DE SCHRYVER G.-M. (2003), *Lexicographers' Dreams in the Electronic-dictionary Age*, in *International Journal of Lexicography*, 16.2: 143-199.
- GAUME B. (2004), “Ballades aléatoires dans les Petits Mondes Lexicaux”, in: *I3, Information Interaction Intelligence*, 4.2 : 1-59 ([w3.univ-tlse2.fr/erss/textes/pagespersos/gaume/resources/I3.impression.5.pdf](http://w3.univ-tlse2.fr/erss/textes/pagespersos/gaume/resources/I3.impression.5.pdf)).
- GROSSMANN F. and TUTIN A. (eds.) (2003), “Les collocations. Analyse et traitement”, in *Travaux et recherches en linguistique appliquée*, série E, n° 1.
- HERBST T. and POPP K. (1999), *The Perfect Learners' Dictionary (?)*, Tübingen.
- RUNDELL M. (1998), “Recent trends in pedagogical lexicography”, in: *International Journal of Lexicography*, 11.4: 315-342.
- VERLINDE S., SELVA T. and BINON J. (2005), “Dictionnaires électroniques et environnement d'apprentissage du lexique”, in *Revue française de linguistique appliquée*, X.2: 19-30.
- VERLINDE S., SELVA T. and BINON J. (2006), “The *Base lexicale du français*: a Multifunctional Online Database for Learners of French”, in Corino E., Marelllo C., Onesti C. (eds.), *Proceedings XII Euralex International Congress. Torino, Italia, September 6th-9th 2006*, Torino, vol. II: 471-481.
- VERLINDE, S., BINON J., OSTYN S. and Bertels A. (to appear), “La Base lexicale du français (BLF): un portail pour l'apprentissage du lexique français”, in *Cahiers de Lexicologie*.



# Linguistic anomaly

Carl Vogel  
Trinity College Dublin

## 1. Introduction

My work in computational linguistics is influenced by the course of my education: a liberal arts undergraduate degree in computing, literature, philosophy and psychology; an MSc by research in artificial intelligence focussed on inheritance reasoning for constraint based syntax; a PhD in cognitive science on models of default reasoning and their relation to human reasoning.

## 2. Outline of the presentation/demo

The direction of research is attuned to cognitive science paradigms, within the field of computational linguistics. Thus, I am interested in formal language theory from the perspective of, for example, theoretical learnability results more than the impact on engineering artifacts, although I am also interested in those. The Irish language spelling checker licensed by the most ubiquitous software company is an example. From my perspective, the practical value was in providing support for daily use of the Irish language in working settings taken for granted for immediate linguistic support and feedback by speakers of the world's major languages, and the theoretical interest was in devising both a space-efficient data structure (think of it as an optimized two pointed trie – something like a rugby ball more than a trie) for a large wordlist, and a time-efficient method of populating it.

This has turned into a focus on linguistic anomaly. However, the cognitive science background entails empirical and theoretical interest, in general.

While I have done work in descriptive syntax within Head-Driven Phrase Structure Grammar (HPSG), particularly on quirky case, I have also worked on parsing for HPSG and formal foundations of HPSG attending to the feature logic. In particular, I am interested in paraconsistent feature logic for a version of HPSG that affords description of degrees of grammaticality. In this context, the related work with and by Jennifer Foster has been useful.

With respect to semantics, like many, I have addressed underspecification – compact representations of ambiguity that model the modest human overload that comes with

ambiguous language, yet the impressive facility humans have for reasoning with only partial resolution of ambiguity. However, I have also considered “overspecification” the process by which accepted linguistic expressions have their senses extended to new meanings and the constraints that exist both theoretically, and in human behavior, with respect to sense extension. This is intimately linked to metaphoricity.

I have also studied study of logics of human reasoning with defaults, chains of statements that express regularities confronted with exceptions. Here there are concerns with the formal properties of the logics themselves, and with the degree to which they serve as adequate models of human reason with generalizations that have exceptions.

Thus, my interests in linguistic anomaly span from orthographical well-formedness, through appropriateness of lexical meaning, formal syntactic description and degrees of grammaticality, to semantic well-formedness and sense extension for representations, and reasoning with incomplete and inconsistent information.

Currently, my funded research is in techniques for text classification, looking particularly at linguistic change over time, towards establishing milestones of normal language development and decline, particularly addressing the Iris Murdoch corpus as a source of data that may reveal features that correlate with progression of Alzheimer’s disease.

The uniting theme in all of these sorts of linguistic anomaly that I study is the tension between linguistic convention and linguistic creativity: ill-formedness versus creativity.

### **3. Issues and challenges**

The practical challenge for work on semantics and reasoning is in advancing beyond spelling checkers and text entailment contests to semantic checkers, tools that can be used to improve drafting of legislation and support human decision making.

The main empirical challenges are those shared with psychologists and psycholinguistics generally, in defining experiments that meaningfully test theoretical claims.

Large empirical challenges are also associated with corpus linguistics: the longitudinal analysis of language change and correlations between milestones of language change and life events depends on a kind of corpus collection that simply does not adequately exist at the present.

### **4. References – Some Representative Publications**

APPEL C. and VOGEL C. (2001), “Investigating Syntax Priming in an E-Mail Tandem Language Learning Environment”, in K. Cameron (ed.), *C.A.L.L. – The Challenge of Change: Research and Practice*, Elm Bank Publications, Exeter: 177-184.

- DEVITT A. and VOGEL C. (2003), "Using Wordnet Hierarchies to Pinpoint Differences in Related Texts", in *Proceedings of EUROLAN 2003: Ontologies, and Information Extraction International Workshop*: 37-44.
- EBERLE K. and VOGEL C. (2000), "Compact Representations of Ambiguous Language", in N. Nicolov and R. Mitkov(eds), *Recent Advances in Natural Language Processing II*, John Benjamins: 191-206.
- FOSTER J. and VOGEL C. (2004). "Parsing Ill-Formed Text Using an Error Grammar", in *Artificial Intelligence Review*, 21: 269-291.
- HEALEY P.G.T., VOGEL C., and Eshghi A. (2007), "Group Dialects in an Online Community", in R. Arnstein and L. Vieu (eds), *DECALOG 2007, The 10th Workshop on the Semantics and Pragmatics of Dialogue*, Università di Trento, May 30 - June 1, 2007: 141-147.
- POPOWICH F. and VOGEL C. (1991b), "A Logic Based Implementation of Head-Driven Phrase Structure Grammar", in C. Brown and G. Koch (eds), *Natural Language Understanding and Logic Programming, III*, Elsevier, North-Holland: 227-246.
- SCHOETER A. and VOGEL C. (eds), (1995), *Nonclassical Feature Systems*, Vol. 10 of *Edinburgh Working Papers in Cognitive Science*, University of Edinburgh, 216 p.
- VAN GIJSEL S. and VOGEL C. (2003), "Inducing a Cline from Corpora of Political Manifestos", in M. Aleksey *et al.* (eds), *International Symposium on Information and Communication Technologies* : 304-310.
- VOGEL C. (2001), Dynamic Semantics for Metaphor, in *Metaphor and Symbol*, 16 (1 & 2): 59-74.
- VOGEL C. (2007), "N-gram Distributions in Texts as Proxy for Textual Fingerprints", in A. Esposito, E. Keller, M. Marinaro and M. Bratanić (eds), *The Fundamentals of Verbal and Non-Verbal Communication and the Biometrical Issue*, IOS Press: 189-194.
- VOGEL C. and POPOWICH F. (1997), "A Parametric Definition for a Family of Inheritance Reasoners", in *New Generation Computing*, 15: 247-292.
- VOGEL C. and VILLADA B. (2000), "Spanish Psychological Predicates", in R. Cann, C. Grover and P. Miller (eds), *Grammatical Interfaces in Head-driven Phrase Structure Grammar*, CSLI Publications, Stanford: 251-66.
- VOGEL C. (1996), "Human Reasoning with Negative Defaults", in D. Gabbay and H.J. Ohlbach (eds), *Practical Reasoning*: 606-621. Lecture Notes in Artificial Intelligence 1085, Springer Verlag, Berlin. *Proceedings of the International Conference on Formal and Applied Practical Reasoning, FAPR'96*, Bonn, June 1996.



# **Evolving approaches to web-supported language learning: From platforms to platform-independent tools**

David Wible  
National Central University

## **1. Introduction**

The task of designing systems and tools that support language learning on the Web is changing due to evolving demands on such technologies from various sources. The design of Learning Activity Management Systems (LAMS) dedicated to language learning, for example, is facing a countervailing trend toward the use of all-purpose platforms such as Blackboard and Moodle. Parallel trends toward the consolidation of systems and content creation include standards specification movements such as SCORM. This overall convergence of content standards and a few all-purpose platforms does not represent an unqualified positive benefit for language pedagogy. The advantages are limited by the unique nature of language learning among learning domains and by the growing availability of digital language tools that ignore both SCORM conformity and portability into larger platforms. Finally, the communicative turn in language pedagogy highlights the need for individualized learning experiences suited to each learner's communicative needs and interests. This sort of individualization is traditionally the forte of ITS (Intelligent Tutoring Systems), yet few of these focus on language learning; those that do are stand-alone systems having virtually no interoperability with other platforms; they treat narrowly-defined dimensions of language learning; and they do not extend or scale up easily or at all.

In this talk, I describe some ongoing work by our team in Taiwan that addresses these current challenges in the design of Web-supported language learning technologies.

## **2. Outline of the presentation/demo**

The first part of the presentation sketches a dedicated language learning platform that has been created and implemented in Taiwan over the past 8 years. IWiLL is the most widely used language learning platform in Taiwan (<http://www.iwillnow.org>). Since 2000, it has been used by more than 190 schools in Taiwan, over 600 teachers, more than 20,000 students. As an online environment, it provides authoring tools for

interactive, multimedia web-based language activities and lessons. It tracks and profiles learner activity and teacher feedback, and automatically stores learners' writing in a dynamic learner corpus (over 3 million words indexed with over 70,000 tokens of teacher feedback).

The second part of the presentation describes a recent collaboration between the IWiLL design team in Taiwan and City University of Hong Kong aimed at modularizing the writing components of IWiLL, creating building blocks or modules from IWiLL's unique language teaching and learning functions so they can be imported into City University's Blackboard platform.

Part three describes a further step in this diffusion away from a dedicated language learning platform toward flexible, modular tools. This step involves browser-based language learning tools which are completely platform independent and accompany learners on their unrestricted navigation of the Web. These ubiquitous tools detect linguistic features in real time on the web pages that the user freely browses and discretely offers these for the selective attention of the learner. This approach is illustrated with a tool designed by our team in Taiwan that focuses on collocations, called Collocator.

### 3. Issues and challenges

The main issue in this research is how to create technologies that have high portability and integrate well in existing Web environments. The most portable of these novel technologies face the additional challenge of being required to perform in real time under noisy, unscripted conditions on the Web. These challenges result from our approach, which contrasts with more widely known stand alone language tools that do not integrate well within existing Web environments or that rely on prescribed content or fragile and baroque learner models that do not scale well.

### 4. References

- IKEDA M., ASHLEY D.K. and CHAN T.-W. (eds), *Intelligent Tutoring Systems: Eighth International Conference, ITS 2006, Lecture Notes in Computer Science 4053*, Springer, Berlin.
- WIBLE D. (in press), "Multiword Expressions and the Digital Turn", in F. Meunier and S. Granger (eds), *Phraseology in Language Learning and Teaching*, John Benjamins, Amsterdam.
- WIBLE, D. (2005), *Language Learning and Language Technology: Toward Foundations for Interdisciplinary Collaboration*, Crane, Taipei.
- WIBLE D., KUO C.-H., CHEN M.C., TSAO N.-L. and HUNG T.-F. (2006), "A Ubiquitous Agent for Unrestricted Vocabulary Learning in Noisy Digital Environments", in *Lecture Notes on Computer Science*, 4053: 503-512.
- WIBLE, D., KUO C.-H., TSAO N.-L. (2004), "Contextualizing Language Learning in the Digital Wild: Tools and a Framework", in *Proceedings of IEEE International Conference on Advanced Learning Technologies*, Joensuu.



WIBLE D., KUO C.-H., TSAO N.-L., HSIU-LING LIN A.L. (2003), "Bootstrapping in a Language Learning Environment", in *Journal of Computer-Assisted Learning*, vol 19 #1: 90-102.



# The SACODEYL project – Corpus exploitation for language learning purposes

Johannes Widmann  
University of Tübingen

## 1. Introduction

SACODEYL is situated in the field of computer-assisted language learning with the help of recent developments in corpus research.

SACODEYL is a project within the SOCRATES-MINERVA initiative whose main aim is to develop an ICT-based system for the assisted compilation and open distribution of European teen talk. This scheme encompasses two groups: group 1, youngsters between 13 and 15 and, group 2, those between 16 and 18.

The main aim of the project is for young Europeans to use corpora for the learning of languages. The pedagogical rationale of the project rests upon notions of autonomous learning and meaningful interaction. SACODEYL users will come into close contact with the real *voices* of peer young Europeans from other countries, their feelings, opinions and speech, without the mediation, otherwise natural and necessary, of third parties such as publishing houses. These peer group voices will make it easier for young people to identify with the language and the contents being taught.

The SACODEYL project aims inter alia at developing a pedagogically-driven search tool for querying corpora in such a way that allows language teachers to access corpora from their teaching perspectives and experiences. The basis of the corpus annotation will be the SACODEYL annotation and corpus enrichment scheme that will be used with the raw transcripts.

## 2. Outline of the presentation/demo

First, I would like to show the design of our annotation tool and the rationale that has been guiding its development. The novel and most important aspect of this annotation tool as opposed to existing tools is the fact that it is based solely on pedagogical criteria.

On the one hand, this includes traditional grammatical and lexical annotation categories that are based on pedagogical grammars. On the other hand, this includes newer categories, such as communicative functions, references to various CEF scales,

information on typical spoken language properties, and information on texture properties. Teachers will be able to modify the corpus annotation and thus they can include their own information in any way they wish. The important thing is that they won't need any advanced computational knowledge. The annotator produces standard TEI-conform XML files without the user having to know about XML. During the demo session, we will look at a beta version of the annotation tool and also at some annotated interviews that have been produced in the project.

Second, we will be able to look at a prototype of the search tool which can be used to retrieve the data and the annotation of the corpora. This online tool will be available to all interested parties free of charge. The tool provides for different viewing perspectives of the corpus. In contrast to many existing query tools, it provides not only a concordance-based viewing perspective but also a section-based viewing perspective where different search options can be viewed in contrast, such as co-occurrences of words and annotated categories within the same section.

To our knowledge, this is one of the first projects where corpus linguistic methodologies have systematically been matched with pedagogic criteria and where pedagogic criteria have been essential in the design of the software. The project explicitly takes into account teachers' needs and perspectives and it aims at offering a low threshold for those teachers who are not very much acquainted with ICT. All the tools are menu-driven. The search tool is an online jsp program requiring no previous installations.

### **3. Issues and challenges**

Given the feedback that we have received so far, it seems as if the SACODEYL project is filling a gap with its approach that has not yet been adequately addressed. We have had quite a number of teaching sites responding positively to our design.

Some of the challenges that lie ahead can be phrased as follows: So far, very little computational intelligence has been included in the software of the project. All the annotation has to be done by hand. On the one hand, this is positive as it is less daunting for the teachers to start off with corpus-based teaching. On the other hand, the project does not yet take advantage of some of the possibilities that computational linguistics seems to offer. This remains a challenge that asks for closer cooperation between pedagogically-interested linguists and technology-interested linguists.

Furthermore, the design of the web-based management and integration of the corpora is also an issue that needs more attention. We have designed an integrated system, but due to a lack of resources it has not yet been possible to fully implement the system.

### **4. References**

BRAUN S. (2005), "From pedagogically relevant corpora to authentic language learning contents", in *ReCALL* 17:1: 47-64.

- BRAUN S. (2006), “ELISA – a pedagogically enriched corpus for language learning purposes”, in S. Braun, K. Kohn and J. Mukherjee (eds), *Corpus Technology and Language Pedagogy: New Resources, New Tools, New Methods*, Peter Lang, Frankfurt/M.
- MUKHERJEE J. (2004), “Bridging the gap between applied corpus linguistics and the reality of English language teaching in Germany”, in U. Connor and Th. Upton (eds), *Applied Corpus Linguistics: A Multidimensional Perspective*, Rodopi, Amsterdam.

