



**HAL**  
open science

## Corpus language input, corpus processes in learning, learner corpus product. Introduction

Alex Boulton, James Thomas

► **To cite this version:**

Alex Boulton, James Thomas. Corpus language input, corpus processes in learning, learner corpus product. Introduction. Alex Boulton; James Thomas. Input, Process and Product: Developments in Teaching and Language Corpora, Masaryk University Press, pp.250, 2012. hal-00683775

**HAL Id: hal-00683775**

**<https://hal.science/hal-00683775>**

Submitted on 12 Jan 2019

**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.

**Alex Boulton & James Thomas. (2012); Corpus language input, corpus processes in learning, learner corpus product. Introduction to J. Thomas & A. Boulton (eds), *Input, Process and Product: Developments in Teaching and Language Corpora*. Brno: Masaryk University Press, p. 7-34.**

Corpus linguistics is essentially concerned with describing language for linguistic research purposes, but language corpora (along with the associated tools and methodologies) have many different affordances and applications. In the field of language teaching, corpus analysis is used to inform the content decisions of what to teach different learner populations in different contexts at different stages of development. This typically includes the application of frequency data in determining the sequence in which linguistic items should be introduced, in identifying key multi-word units and a wide range of lexico-semantic patterns, and in predicting areas of potential difficulty from learner corpora. This essentially indirect approach (Römer 2011) to corpus data is taken by syllabus designers, materials writers, lexicographers and testers, though the results may be entirely invisible to the end user (McCarthy 2004). However, teachers can also make use of corpora to answer their own questions about language, to test grammar ‘rules’ against real data, to find examples and help create materials for teaching and testing, among others. Learner involvement need not be limited to teacher mediated uses, but can involve direct hands-on consultation, either for language learning or as a reference resource. This is commonly associated with the work of Tim Johns<sup>1</sup> in what he called data-driven learning (DDL), an approach which conflates the roles of learners and researchers and sees them deriving their own answers from direct contact with the data (e.g. Johns & King 1991). The approach is essentially constructivist, providing an authentic way of tackling lexico-grammar in particular (Thomas 2006) in contrast to most decontextualised and relatively ‘artificial’ vocabulary learning techniques – assuming any strategies are taught at all.

The papers selected for inclusion in this volume derive from presentations given at TaLC9 in Brno in 2010, 16 years after the first TaLC conference was held in Lancaster in 1994. Looking through the list of over 150 papers published from almost two decades of TaLC conferences (see Appendix A), an evolutionary trajectory emerges: while many of the early issues are still relevant today, other have opened up in various ways, and this volume includes some papers that cover entirely new ground. TaLC is thus no longer in its infancy – but neither has it reached full maturity. It has gone beyond the initial idea of concordancing by advanced adult L2 students for lexico-grammar (e.g. Tribble & Jones 1997), to being employed in an ever-expanding array of linguistic fields from discourse analysis (e.g. Charles 2007) to literary studies (e.g. Kettemann & Marko 2004, 2011) to translation (e.g. Kübler 2003), at lower levels (e.g. Cobb et al. 2001), in schools (e.g. Sun & Wang 2003) and even for primary schools for L1 (e.g. Sealey & Thompson 2007). The early research enthusiasm is as strong as ever and is constantly passed on to generations of new researchers, but given the development of new types of corpora, of more sophisticated software and of computer technology in general, there can be no certainties about what directions it is likely to take, nor how it may eventually earn its keep in regular classroom practice. Despite the considerable technological advances and numerous publications in the spheres of language education, it

---

<sup>1</sup> 1936-2009. See the obituary by Scott (2009).

is frequently remarked that TaLC remains marginal to mainstream language teaching (e.g. Chambers et al. 2011).

One probable cause of this lack of uptake in mainstream language education is that TaLC, at least in popular perception, remains stubbornly the province of researchers rather than teachers (Mukherjee 2004), a gap that desperately needs to be bridged (cf. McCarthy 2008). Worse, corpus work is seen as an ivory tower activity, generating a notable lack of empirical classroom research (e.g. Johansson 2009; Yoon 2011). However, a growing body of studies do attempt to evaluate some aspect of corpus use in real classroom contexts – 93 separate studies to date, according to a current survey by Boulton (2010). These are tremendously varied in design, underlining the flexibility of approaches to corpus use for a variety of different learner needs in very different conditions; as Breyer (2006: 162) has pointed out, corpus activities are “limited only by the imagination of the user.” But however corpora are introduced, the overwhelming conclusion is that learners can use them effectively for many different purposes, are receptive to the approach and see the relevance to their own needs, and can use them successfully both as a learning tool and as a reference resource, particularly for writing, revision, error-correction and translation.

The TaLC conference series combines, as its name suggests, teaching and language corpora. But crucially, teaching is the first of the two terms, and this is reflected in the structure of the present volume, with the first two sections looking at how corpora can be used as *input* for language learning. Section One opens with a paper by Ana Frankenberg-Garcia, who asks why corpus use is not more widespread among the language teaching community, and provides a number of suggestions for how corpora can be integrated into everyday language classes. For her, the crucial issue is not what teachers and learners can do with corpora, but what corpora can do for teachers and learners. The remaining chapters in this section explore some of the potential for corpus use in language teaching. Patrick Hanks combines prototype theory and corpus linguistics to show how pattern analysis can lead to a radically different approach to language and linguistics, in the process transforming dictionaries and other reference resources for language teachers and learners. The result is firmly rooted in actual language use, integrating focus on form and on meaning into a fundamentally innovative tool for these end users. Teachers and learners can also exploit corpora as a reference resource, as discussed largely in Section 2, but a number of initial considerations in developing corpora and software are reported in the next two papers in this section. Shozo Yokoyama, Chizuko Suzuki, Seisuke Yasunami and Naoko Kawakita describe the construction of a corpus of academic research articles in medicine, which they analyse for different types of verbs. It is argued that learners can benefit from the resulting insights in terms of frequency, keyness, collocates and distributions over different IMRAD sections, which they can discover using the dedicated corpus interface outlined in the paper. Ute Römer also describes a specialised corpus and interface, but here compiled from high-scoring essays mainly by native speakers who are still learning their own discipline. The Michigan Corpus of Upper-level Student Papers (MICUSP) is thus pedagogically relevant to EAP learner / apprentice writers: teachers can use it to inform their teaching, and learners can explore it in a DDL approach to academic writing through a simple on-line interface, as the paper reports. MICUSP is the written counterpart to MICASE (the Michigan Corpus of Academic Spoken English), and though corpora of spoken language are more difficult to compile than those of

written language, they are of great importance in developing the teaching of oral skills. To this end, Stefanie Dose shows that a corpus of TV transcripts can be tremendously valuable for pedagogical purposes, demonstrating that the language is remarkably similar to unscripted speech. TV series can provide a corpus that learners can relate to or 'authenticate' (cf. Widdowson 2000), and allow work on individual written or multimedia extracts for a variety of activities – a "pedagogically relevant" corpus in Braun's (2005) terms. While we can certainly subvert linguistic corpora for language teaching, this inevitably involves a certain amount of "rethinking" (Burnard & McEnery 2000).

These introductory chapters derive from the contributors' many years of experience in using corpus data either directly or indirectly for language learning – they are far from ivory tower expositions divorced from reality. Section Two makes the connection between corpus and classroom more explicit: all of the contributions report on actual applications and evaluate outcomes, attitudes and behaviours of learners faced with corpora and associated tools – the processes involved in using corpora in language teaching and learning.

A recurring question is how corpus work can be successfully integrated to normal classroom practice, as highlighted in the paper by Monika Geist<sup>2</sup> and Angela Hahn. Their results are encouraging insofar as their learners are clearly able to use the general British National Corpus (BNC) for specific ends with some success, even though some of them lacked the necessary motivation to invest time and effort in corpus activities which were not graded and which the learners were unable to relate to their regular classes. It is common practice to introduce corpus activities as an add-on, going against the precept of constructive alignment (e.g. Biggs 1996). But DDL can be introduced as 'ordinary' practice as demonstrated in the study by Henry Tyne, who shows that it is perfectly compatible with standard teaching techniques – including at the level of text. The teachers in his study report that the DDL techniques involved are of immediate benefit in their daily teaching, and may even provide a way in to more usual DDL activities later on. Another option is for the teacher to mediate the corpus data and use only printed materials, thus eliminating the 'obstacle' of the computer in DDL. Alex Boulton reports on using DDL with and without a computer, finding that each approach has its own advantages in terms of learning outcomes and appeals to different learners. In a similar vein, Kiyomi Chujo and Kathryn Oghigian find that optimal results may be obtained from a combination of paper-based and computer-based DDL, here in terms of feedback and learning outcomes for vocabulary and grammar. Examples such as these show that corpora can be easily and efficiently exploited by learners even without extensive training in the associated tools. This is confirmed in the following paper by Klára Osolobě and Pavlína Vališová, where learners of Czech managed to conduct simple queries and obtain meaningful results with a minimum of training. Even the seemingly complex work with lexical bundles reported by Andreas Eriksson was conducted over only two workshop sessions, suggesting that focusing on specific tasks in relevant specialist fields can make corpus work more relevant and motivating and thus more accessible.

---

<sup>2</sup> Monika Geist originally contributed to this paper as Monika Formánková.

These first two sections show that corpus use is no longer the sole preserve of the “particular type of student” typical of early DDL work – “adult: well-motivated, a sophisticated learner with experience of research methods in his subject area with particular needs... in a particular learning/teaching situation” (Johns 1986: 161). This evolution is perhaps inevitable with the increasing availability of a variety of corpora and more user-friendly software, appropriate even for secondary school students as exemplified in the studies by Geist and Hahn as well as by Tyne (where the teachers are also regular teachers and not researchers). Though it is true that many of the studies here do involve undergraduates, most are students who are not majoring in languages, often with low levels of motivation, little sophistication in language learning, and relatively low levels of proficiency – pre-intermediate in Boulton, beginners in Chujo and Oghigian.

While English is perhaps inevitably the most common target language, Tyne’s students are learning Spanish, Osolobě and Vališová’s learning Czech (one cohort even consists of native speakers), underscoring the flexibility of corpus-based activities even for languages which are quite different from English in terms of morphological complexity and syntax. The types of data used also vary widely, from four million words of general English in Geist and Hahn to the level of individual text in Tyne; from student papers in Römer to expert writing in Eriksson and Yokoyama et al.; parallel corpora in Chujo and Oghigian; spoken data in Dose, and so on. The tasks and types of analysis are correspondingly varied, from the very simple lexical level for younger learners in Geist and Hahn to lexical bundles in Eriksson and phraseology in Römer. The overall picture which emerges is that corpora and DDL hold something for everyone: there is no ‘best’ corpus for all purposes and no exclusive ‘right’ way to exploit corpora: pedagogical relevance and appropriateness in each specific case is paramount (Flowerdew 2009).

Sections 3 and 4 move on to learner corpora, i.e. corpora compiled from the spoken or written *output* of learners, which can be quantified and analysed in the same way as corpora consisting of native or expert texts (Leńko-Szymańska 2008). The results serve many purposes as can be seen from the wide variety of issues covered here, reflecting the burgeoning field of learner corpus research spanning the last 20 years (cf. Granger 2009). As with corpora of native speaker or expert texts, learner corpora can be used in a data-driven learning approach (Granger & Tribble 2006) where learners analyse corpora comprising texts of their own language output or those of others (Seidlhofer 2000). They are also valuable in the automatic detection of errors and the automatic correction and scoring of student writing. They can be used to inform materials, resources and practices as well as testing and assessment tools. They can improve our knowledge of the processes involved in language acquisition and interlanguage development, and allow us to relate particular features to different levels of proficiency. In the classroom, they are a resource for systematically raising teachers’ awareness of their own learners’ specific problems. And on the positive side of the coin, the successful use of specific features of student output can be observed and used as models of good practice.

But probably the most frequent approach, and the one that launches Section Three, is the comparison of learner and native corpora, usually with a focus on ‘errors’ – including the under- and overuse of various linguistic features. Corpus linguistics allows rigorous analysis

or learner output for systematic detection and analysis of areas of difficulty where previous attempts could rely on little more than a hunch based on personal experience or intuition –; it is therefore unsurprising that contrastive analysis has made something of a comeback in recent years. Several papers here thus attribute different error types directly to the learner’s mother tongue (L1), potentially an argument for a return to the use of materials produced with the specific L1 in mind and against the use of generic textbooks produced by international publishers for global distribution.

Marina Mattheoudakis and Anna-Maria Hatzitheodorou compare learner writing against native texts for collocates of delexical or ‘light’ verbs. Their analysis suggests that transparency and the existence of comparable collocates in the L1 are major factors in predicting erroneous as well as over- and underused collocates; without them, learners have little choice but to rely on Sinclair’s (1991: 109ff) “open-choice principle” rather than his “idiom principle”. As such items tend to lack salience, training is needed in noticing. This is the case for many spoken features too, as shown in the paper by Sandra Götz who finds that even advanced learners tend to speak less (in terms of words per minute or length of turn) than native speakers, and exhibit greater use of unfilled pauses and other hesitation phenomena along with more limited use of discourse markers. A final paper comparing learner and native corpora also looks at discourse markers in speech: Jiajia Xu, Mark Morgan and John McKenny highlight the need for intuition in complementing automatic extraction of semantically relevant n-grams. Differences are again attributed largely to L1 transfer, with overuse in particular being linked to a more limited repertoire of connectors due in part to decontextualised overteaching of specific items. A similar point is made by Svetla Rogatcheva, who contrasts required and optional contexts for different verb aspects in the present and past, showing that Bulgarian learners have more difficulty with the English progressive, German learners with the perfect. These problems can be linked not only to the L1, but also again to overteaching which might deter learners from using items perceived as problematic. Most of these papers are based on existing learner corpora, but Sylwia Twardo shows that it is possible to create even a fairly large (300,000-word) PoS-tagged learner corpus from scratch. She takes up a theme mentioned by Rogatcheva and Xu et al., namely the difficulties involved in dealing with automatic error-detection. These are most visible in the form of ‘non-words’ arising from spelling or morpheme errors, which occur fairly predictably across different levels of proficiency.

Such contrastive analyses are certainly useful, but the authors do not claim that every difference between native and non-native use is an error to be eradicated at the earliest opportunity: there is often a good reason underlying interlanguage differences (Aston 2008). For example, the presence or overuse of some features (e.g. full forms instead of contractions, overuse of connectors or temporal markers) may increase communicative effectiveness if they in fact compensate for other difficulties (e.g. mastery of pronunciation, deixis or tenses respectively). Similarly, the absence or underuse of particular items (e.g. complex sentence structures or phrasal verbs) may also be communicatively more effective at early stages of development (cf. Larsen-Freeman & Cameron 2008). Finally, learners may even be more effective than monolingual native speakers in intercultural contexts where they may, for example, use fewer idioms or opaque expressions, and be more direct in speech acts such as disagreeing or asking for help. While it is important to note such

differences, for all these reasons care should be taken to distinguish features that significantly impede communication, those that have little if any effect, and those that may actually be advantageous (cf. Seidlhofer 2011). The point being made here is that the value of learner corpora goes beyond mere error analysis, and it is as important to see what learners *can* do as what they can't – all, of course, for different learners in different conditions at different stages of development (cf. the earlier discussion of MICUSP by Römer).

These are some of the issues taken up in the final section of highly innovative papers, beginning with the article by Susanne Kämmerer: although she also discusses errors in a series of studies, this is crucially from the learner's perspective. Three years after the compilation of the corpus, the original German contributors were able to detect their own errors in only 30% of cases; however, they were able to correct almost all errors once they were pointed out and to explain most, attributing them overwhelmingly to L1 interference or 'stupid mistakes'. Such insights are important, as the inevitable question is what a teacher should do with errors once they have been detected. M. Trevor Shanklin addresses this issue in considering how automatically generated feedback from oral exams should be useful not just to test-designers and examiners but also to test-takers. This is the aim of the corpus in the Contrastive Analysis Screening Tool (CAST): basic information such as type/token ratio and mean length of utterance are discussed in relation to proficiency, as are more specific features such as the appropriate use of tenses and subordination. While much of this still focuses on errors, the intention is for the corpus to further serve as an indicator of what successful learners can actually do at different levels, an assumption underpinning the English Vocabulary Profile lists analysed in the final paper by Yukio Tono. The underlying idea of the English Profile project (now with its own journal) is to provide detailed descriptions of what learners of English show they can do at different levels rather than identifying what they get wrong (i.e. what they *should* know). This laudable aim is inevitably fraught with difficulties, as Tono's analysis reveals: in particular, the procedures for deriving the lists from the very large Cambridge native and learner (exam) corpora are not entirely transparent, and it is difficult to attribute different levels to the different senses and uses of individual items. The problems are similar in this respect to the sequencing of dictionary entries, but it is argued that particular attention needs to be paid to receptive and productive uses.

Most of the papers in these sections on learner corpora use a published corpus, especially one of those made available at the Centre for English Corpus Linguistics (CECL) at the Université Catholique de Louvain<sup>3</sup>, namely the International Corpus of Learner English (ICLE) and the Louvain International Database of Spoken English Interlanguage (LINDSEI). The former consists of written texts in the form of argumentative essays, the second of traditional oral exam-style questions. One advantage of this suite of corpora is that it is possible to focus on a sub-corpus of learners according to their L1: CECL sub-corpora from Bulgarian, French, German, Greek and Spanish learners all feature in the papers here, along with L1 Chinese and Polish from other sources. Only Shanklin and Tono compound learner corpora from speakers of different L1s, but for very explicit reasons: in the former, to

---

<sup>3</sup> See <http://www.uclouvain.be/en-cecl.html>, accessed 20/11/11.

produce tools that can be used for different target languages; in the latter to explore a generic, non-language specific resource from a major publisher.

ICLE and LINDSEI can each be compared against an equivalent native speaker corpus also produced by the CECL: the Louvain Corpus of Native Speaker English Essays (LOCNESS), and the Louvain Corpus of Native Speaker Conversation (LOCNEC) respectively – the former used in Mattheoudakis and Hatzitheodorou, the latter in Götz. The learner corpora are undoubtedly ‘authentic’ even though the data are gathered in highly controlled conditions, as the contexts reflect ‘typical’ learner communicative contexts – participating in written and oral exams (cf. Mendikoetxea et al. 2010: 183). While the native speaker corpora might be considered less authentic (or at least, less ecological, as native speakers do not necessarily participate in similar types of exams), it clearly makes sense to compare learner language against native language gathered in comparable situations. However, other corpora such as MICASE or the BNC are for many purposes sufficiently comparable (as here in Xu et al.).

TaLC, then, is maturing nicely. Kudos must of course go to the visionary pilgrim fathers who made the connection between esoteric linguistic research and the overwhelmingly practical concerns of language teaching and learning, but the ever-expanding CV of TaLC-related publications<sup>4</sup> bears testament to growing research interest around the world. And not just research: the various corpora at Brigham Young University are accessed by over 80,000 individual users each month; of these, only 15% declare their main interest in corpora as being for research purposes (in linguistics, sociology, cultural studies, literature and politics); 28% for professional uses (translators, writers, lexicographers and testers). 15% are teachers (native and non-native), but the largest group by far consists of language learners at 42%.<sup>5</sup> This augurs well for further developments relating teaching and language corpora, an area to which this volume makes its own contribution.

The present volume would not have been possible without the input of certain individuals and organisations. First among these is the TaLC organising committee who blind-reviewed the papers prior to the Brno conference (2010) as well as all full submissions to this volume: Guy Aston, Lou Burnard, Lynn Flowerdew, Bernhard Kettmann, Natalie Kübler, Agnieszka Leńko-Szymańska, Ute Römer and Christopher Tribble. We are also enormously grateful to Marek Procházka, a doctoral student in the Faculty of Arts at Masaryk University, for his typesetting of the whole book.

## References

- Aston, G. 2008. It’s only human... In A. Martelli & V. Pulcini (eds), *Investigating English with Corpora: Studies in honour of Maria Teresa Prat*. Monza: Polimetrica, p. 343-354.
- Biggs, J. 1996. Enhancing teaching through constructive alignment. *Higher Education* 32: 347-364.
- Boulton, A. 2010. Learning outcomes from corpus consultation. In M. Moreno Jaén, F. Serrano Valverde & M. Calzada Pérez (eds), *Exploring New Paths in Language Pedagogy: Lexis and corpus-based language teaching*. London: Equinox, p. 129-144. Updated

---

<sup>4</sup> See the bibliographical database at CorpusCALL, the corpus-related Eurocall SIG: <http://corpuscall.eu/>, accessed 24/11/11.

<sup>5</sup> <http://corpus.byu.edu/>, accessed 24/11/11. Figures kindly provided by Mark Davies, personal communication.

- supplement (description of 93 empirical DDL studies) at CorpusCALL:  
<http://corpuscall.eu/course/view.php?id=5#sectionblock-2>, accessed 15/09/11.
- Braun, S. 2005. From pedagogically relevant corpora to authentic language learning contents. *ReCALL* 17(1): 47-64.
- Breyer, Y. 2006. My Concordancer: Tailor-made software for language learners and teachers. In S. Braun, K. Kohn & J. Mukherjee (eds), *Corpus Technology and Language Pedagogy: New resources, new tools, new methods*. Frankfurt: Peter Lang, p. 157-176.
- Burnard, L. & T. McEnery (eds). 2000. *Rethinking Language Pedagogy from a Corpus Perspective*. Frankfurt: Peter Lang.
- Chambers, A., F. Farr & S. O’Riordan. 2011. Language teachers with corpora in mind: From starting steps to walking tall. *Language Learning Journal* 39(1): 85-104.
- Charles, M. 2007. Reconciling top-down and bottom-up approaches to graduate writing: Using a corpus to teach rhetorical functions. *Journal of English for Academic Purposes* 6(4): 289-302.
- Cobb, T., C. Greaves & M. Horst. 2001. Peut-on augmenter le rythme d’acquisition lexicale par la lecture? Une expérience de lecture en français appuyée sur une série de ressources en ligne. In P. Raymond & C. Cornaire (eds), *Regards sur la Didactique des Langues Secondes*. Montréal: Editions Logique, p. 133-153. [Translation: Can the rate of lexical acquisition from reading be increased? An experiment in reading French with a suite of on-line resources.] <http://www.lextutor.ca/cv/>, both accessed 01/12/11.
- English Profile Journal*. [on line] <http://journals.cambridge.org/action/displayJournal?jid=EPJ>, accessed 24/11/11.
- Flowerdew, L. 2009. Applying corpus linguistics to pedagogy: A critical evaluation. *International Journal of Corpus Linguistics* 14(3): 393-417.
- Granger, S. 2009. The contribution of learner corpora to second language acquisition and foreign language teaching: A critical evaluation. In K. Aijmer (ed.), *Corpora and Language Teaching*. Amsterdam: John Benjamins, p. 13-32.
- Granger, S. & C. Tribble. 1998. Learner corpus data in the foreign language classroom: Form-focused instruction and data-driven learning. In S. Granger (ed.), *Learner English on Computer*. London: Longman, p. 199-209.
- Johansson, S. 2009. Some thoughts on corpora and second-language acquisition. In K. Aijmer (ed.), *Corpora and Language Teaching*. Amsterdam: John Benjamins, p. 33-44.
- Johns, T. 1986. Micro-Concord: A language learner’s research tool. *System* 14(2): 151-162.
- Johns, T. & P. King (eds.) 1991. *Classroom Concordancing*. *English Language Research Journal* 4.
- Kettemann, B. & G. Marko. 2004. Can the L in TaLC stand for literature? In G. Aston, S. Bernardini & D. Stewart (eds), *Corpora and Language Learners*. Amsterdam: John Benjamins, p. 169-193.
- Kettemann, B. & G. Marko. 2011. Data-driving critical discourse analysis. In N. Kübler (ed.), *Corpora, Language, Teaching, and Resources: From theory to practice*. Bern: Peter Lang, p. 19-48.
- Kübler, N. 2003. Corpora and LSP translation. In F. Zanettin, S. Bernardini & D. Stewart (eds), *Corpora in Translator Education*. Manchester: St Jerome Publishing, p. 25-42.
- Larsen-Freeman, D. & L. Cameron. 2008. *Complex Systems and Applied Linguistics*. Oxford: Oxford University Press.

- Leńko-Szymańska, A. 2008. Non-native or non-expert? The use of connectors in native and foreign language learners' texts. *Acquisition et Interaction en Langue Etrangère* 27: 99-108.
- McCarthy, M. 2004. *Touchstone: From corpus to coursebook*. Cambridge: Cambridge University Press.  
[http://www.cambridge.org/other\\_files/downloads/esl/booklets/McCarthy-Touchstone-Corpus.pdf](http://www.cambridge.org/other_files/downloads/esl/booklets/McCarthy-Touchstone-Corpus.pdf), accessed 07/10/11.
- McCarthy, M. 2008. Accessing and interpreting corpus information in the teacher education context. *Language Teaching* 41(4): 563-574.
- Mendikoetxea, A., S. Bielsa & P. Rollinson. 2010. Focus on errors: Learner corpora as pedagogical tools. In M-C. Campoy, B. Bellés-Fortuño & M-L. Gea-Valor (eds), *Corpus-Based Approaches to English Language Teaching*. London: Continuum, p. 180-194.
- Mukherjee, J. 2004. Bridging the gap between applied corpus linguistics and the reality of English language teaching in Germany. In U. Connor & T. Upton (eds), *Applied Corpus Linguistics: A multidimensional perspective*. Amsterdam: Rodopi, p. 239-250.
- Römer, U. 2011. Corpus research applications in second language teaching. *Annual Review of Applied Linguistics* 31: 205-225.
- Scott, M. 2009. In memory of Tim Johns. *International Journal of Corpus Linguistics* 14(3): 271-274.
- Sealey, A. & P. Thompson. 2007. Corpus, concordance, classification: Young learners in the L1 classroom. *Language Awareness* 16(3): 208-223.
- Seidlhofer, B. 2000. Operationalizing intertextuality: Using learner corpora for learning. In L. Burnard & T. McEnery (eds), *Rethinking Language Pedagogy from a Corpus Perspective*. Frankfurt: Peter Lang, p. 207-223.
- Seidlhofer, B. 2011. *Understanding English as a Lingua Franca*. Oxford: Oxford University Press.
- Sinclair, J. 1991. *Corpus, Concordance, Collocation*. Oxford: Oxford University Press.
- Sun, Y-C. & L-Y. Wang. 2003. Concordancers in the EFL classroom: Cognitive approaches and collocation difficulty. *Computer Assisted Language Learning* 16(1): 83-94.
- Thomas, J. 2006. Using corpora in language teaching and learning. *Teaching English with Technology* 6(1): 44-56.  
<http://www.tewtjournal.org/VOL%206/ISSUE%201/VOL%206%20ISSUE%201%20COMPLETE.pdf>, accessed 24/10/11.
- Tribble, C. & G. Jones. 1997. *Concordances in the Classroom*. 2nd edition. Houston: Athelstan.
- Widdowson, H. 2000. On the limitations of linguistics applied. *Applied Linguistics* 21(1): 3-25.
- Yoon, C. 2011. Concordancing in L2 writing class: An overview of research and issues. *Journal of English for Academic Purposes* 10: 130-139.

## Appendix A: Previous TaLC conferences and publications, 1994–2008.

### TaLC 8. Lisbon, Portugal: ISLA – Lisboa. 3<sup>rd</sup> – 6<sup>th</sup> July 2008.

Ana Frankenberg-Garcia, Lynn Flowerdew & Guy Aston (eds). 2011. *New Trends in Corpora and Language Learning*. London: Continuum.

1. Burnard, L. Preface. xv-xviii.
2. Frankenberg-Garcia, A. Introduction. xxviii-xxxiv.  
Flowerdew, L.  
Aston, G.
3. Tono, Y. TaLC in action: Recent innovations in corpus-based English language teaching in Japan. 3-25.
4. Charles, M. Using hands-on concordancing to teach rhetorical functions: Evaluation and implications for EAP. 26-43.
5. Kettemann, B. Tracing the emo side of life: Using a corpus of an alternative youth culture discourse to teach cultural studies. 44-61.
6. Kübler, N. Working with corpora for translation teaching in a French-speaking setting. 62-80.
7. Kaszubski, P. IFAConc: A pedagogic tool for online concordancing with EFL/EAP learners. 81-104.
8. Liu, A. A corpus-based approach to automatic feedback for learners' miscollocations. 107-120.  
Wible, D.  
Tsao, N-L.
9. Cocchetta, F. Multimodal functional-notional concordancing. 121-138.
10. Curado Fuentes, A. Academic corpus integration in MT and application to LSP teaching. 139-152.
11. Warren, M. Using corpora in the learning and teaching of phraseological variation. 153-166.
12. Widmann, J. The SACODEYL search tool: Exploiting corpora for language learning purposes. 167-178.  
Kohn, K.  
Ziai, R.
13. Osborne, J. Oral learner corpora and the assessment of fluency in the Common European Framework. 181-197.
14. De Cock, S. Preferred patterns of use of positive and negative evaluative adjectives in native and learner speech: An ELT perspective. 198-212.
15. Nesi, H. BAWE: An introduction to a new resource. 213-228.
16. Hatzitheodorou, A-M. The impact of culture on the use of stance exponents as persuasive devices: The case of GRICLE and English native speaker corpora. 229-246.  
Mattheoudakis, M.
17. McKenny, J. Polishing papers for publication: Palimpsests or procrustean beds? 247-262.  
Bennett, K.

### TaLC 7. Paris, France: University Paris 7 Denis Diderot / BNF. 1<sup>st</sup> – 4<sup>th</sup> July 2006.

Natalie Kübler (ed.). 2011. *Corpora, Language, Teaching, and Resources: From theory to practice*. Bern: Peter Lang.

1. Kübler, N. Introduction. 9-15.
2. Kettemann, B. Data-driving critical discourse analysis. 19-48.  
Marko, G.
3. Philip, G. '...and I dropped my jaw with fear': The role of corpora in teaching phraseology. 49-68.
4. Boulton, A. Bringing corpora to the masses: Free and easy tools for interdisciplinary language studies. 69-95.
5. Chambers, A. Language learning as discourse analysis: Playing games in a corpus of French journalistic discourse. 97-112.
6. Charles, M. Corpus evidence for teaching adverbial connectors of Contrast: 'However', 'yet', 'rather', 'instead' and 'in contrast'. 113-131.
7. Van Rij-Heyligers, J. Breaking the chains of rhetorics in academia: Corpus-based research as tool for

- transformation in discourse? 133-154.
8. Krausse, S. Semantic preference and semantic prosody in the specialist language class. 155-164.
  9. Schmied, J. Teaching and learning contrastive linguistics using an EU translation corpus with English, German, French and Spanish. 165-181.
  10. Jimenez-Caycedo, J. Gebhard, M. 'Expert-like' elementary narratives: A genre- and corpus-based study of L2 writing development. 185-198.
  11. Diez Bedmar, M. Casas Pedrosa, A. The use of prepositions by Spanish learners of English at University level: A longitudinal analysis. 199-218.
  12. Castagnoli, S. Ciobanu, D. Kunz, K. Kübler N. Volanschi, A. Designing a learner translator corpus for training purposes. 221-247.
  13. Pecman, M. How awareness of lexical combinatorion can improve second language learning: A model for analysing collocations in scientific discourse. 249-261.
  14. Tsaknaki, O. Recognizing proverbs: A method and its applications. 263-272.
  15. Martin, P. A language teaching software program using spontaneous speech corpora. 273-283.
  16. Kraif, O. Tutin, A. Using a bilingual annotated corpus as an academic writing aid: An application for EFL users. 285-297.
  17. Chiari, I. Teaching language variation using Italian corpora. 301-322.
  18. Williams, G. The learner's dictionary and the sciences: Mismatch or no match? 323-340.

**TaLC 6. Granada, Spain: University of Granada. 6<sup>th</sup> – 9<sup>th</sup> July 2004.**

Encarnación Hidalgo, Luis Quereda & Juan Santana (eds). 2007. *Corpora in the Foreign Language Classroom*. Amsterdam: Rodopi.

1. Hidalgo, E. Quereda, L. Santana, J. Foreword. ix-xiv.
2. Chambers, A. Popularising corpus consultation by language learners and teachers. 3-16.
3. Johansson, S. Using corpora: From learning to research. 17-28.
4. Braun, S. Designing and exploiting small multimedia corpora for autonomous learning and teaching. 31-46.
5. Chujo, K. Utiyama, M. Nishigaki, C. Towards building a usable corpus collection for the ELT classroom. 47-69.
6. Lam, P. A corpus-driven lexico-grammatical analysis of English tourism industry texts and the study of its pedagogic implications in English for specific purposes. 71-89.
7. Guo, X. Errors or partial acquisition: A case study of a young English learner's interlanguage. 91-104.
8. Van Rij-Heyligers, J. To weep perilously or W.EAP critically: The case for a corpus-based critical EAP. 105-118.
9. Meunier, F. Gouverneur, C. The treatment of phraseology in ELT textbooks. 119-139.
10. Perez Basanta, C. Rodriguez Martin, M. The application of data-driven learning to a small-scale corpus: Using film transcripts for teaching conversational skills. 141-158.
11. Coffey, S. Investigating restricted semantic sets in a large general corpus: Learning activities for students of English as a foreign language. 161-173.
12. Gesuato, S. How (dis)similar? Telling the difference between near-synonyms in a foreign language. 175-190.
13. Minugh, D. George Bush and the last crusade, or the fight for truth, justice and the American way. 191-205.

14. Papp, S. Inductive learning and self-correction with the use of learner and reference corpora. 207-220.
15. Olivier, N.  
Brems, L.  
Davidse, K.  
Speelman, D.  
Cuyckens, H. Pattern-learning and pattern-description: An integrated approach to proficiency and research for students of English. 221-235.
16. Lavid, J. Contrastive patterns of mental transitivity in English and Spanish: A student-centred corpus-based study. 237-252.
17. Leńko-Szymańska, A. Past progressive or simple past? The acquisition of progressive aspect by Polish advanced learners of English. 253-266.
18. Cresswell, A. Getting to 'know' connectors? Evaluating data-driven learning in a writing skills course. 267-287.
19. Tribble, C. Managing relationships in professional writing. 289-308.
20. Curado Fuentes, A. A corpus-based assessment of reading comprehension in English. 309-326.
21. Louw, B. Truth, literary worlds and devices as collocation. 329-362.

#### **TaLC 5. Bertinoro, Italy: University of Bologna. 27<sup>th</sup> – 30<sup>th</sup> July 2002.**

Guy Aston, Sylvia Bernardini & Dominic Stewart (eds). 2004. *Corpora and Language Learners*. Amsterdam: John Benjamins.

1. Stewart, D.  
Bernardini, S.  
Aston, G. Ten years of TaLC. 1-18.
2. Hoey, M. The textual priming of lexis. 21-41.
3. Tono, Y. Multiple comparisons of IL, L1 and TL corpora: The case of L2 acquisition of verb subcategorization patterns by Japanese learners of English. 45-66.
4. Borin, L.  
Prütz, K. New wine in old skins? A corpus investigation of L1 syntactic transfer in learner language. 67-87.
5. Leńko-Szymańska, A. Demonstratives as anaphora markers in advanced learners' English. 89-107.
6. Nesselhauf, N. How learner corpus analysis can contribute to language teaching: A study of support verb constructions. 109-124.
7. Flowerdew, L. The problem-solution pattern in apprentice vs. professional technical writing: An application of appraisal theory. 125-135.
8. Chipere, N.  
Malvern, D.  
Richards, B. Using a corpus of children's writing to test a solution to the sample size problem affecting type-token ratios. 137-147.
9. Römer, U. Comparing real and ideal language learner input: The use of an EFL textbook corpus in corpus linguistics and language teaching. 151-168.
10. Kettemann, B.  
Marko, G. Can the L in TaLC stand for literature? 169-193.
11. Mauranen, A. Speech corpora in the classroom. 195-211.
12. Frankenberg-Garcia, A. Lost in parallel concordances. 213-229.
13. Sripicharn, P. Examining native speakers' and learners' investigation of the same concordance data and its implications for classroom concordancing with ELF learners. 233-245.
14. Pérez-Paredes, P.  
Cantos-Gomez, P. Some lessons students learn: Self-discovery and corpora. 247-257.
15. Davies, M. Student use of large, annotated corpora to analyze syntactic variation. 259-269.
16. Fletcher, W. Facilitating the compilation and dissemination of ad-hoc web corpora. 273-300.

#### **TaLC 4. Graz, Austria: University of Graz. 19<sup>th</sup> – 22<sup>nd</sup> July 2000.**

Bernhard Kettemann & Georg Marko (eds). 2002. *Teaching and Learning by Doing Corpus Analysis*. Amsterdam: Rodopi.

1. McEnery, T. TaLC 4: Where are we going? General aspects of corpus linguistics. 3-5.
2. Aston, G. The learner as corpus designer. 9-25.
3. Renouf, A. The time dimension in modern English corpus linguistics. 27-41.
4. Scott, M. Picturing the key words of a very large corpus and their lexical upshots, or getting at the Guardian's view of the world. 43-50.
5. Burnard, L. Where did we go wrong? A retrospective look at the British National Corpus. 51-70.
6. Coxhead, A. The academic word list: A corpus-based word list for academic purposes. 73-89.
7. Mindt, D. A corpus-based grammar for ELT: Data-driven learning. 91-104.
8. Johns, T. Data-driven learning: The perpetual challenge. 107-117.
9. Mair, C. Empowering non-native speakers: The hidden surplus value of corpora in continental English departments. 119-130.
10. Lorenz, G. Language corpora rock the base: On standard English grammar, perfective aspect and seemingly adverse corpus evidence. 131-145.
11. Wible, D. Toward automating a personalized concordancer for data-driven learning: A lexical difficulty filter for language learners. 147-154.
- Chien, F.
- Kuo, C-H.
- Wang, C.
12. Kirk, J. Teaching critical skills in corpus linguistics using the BNC. 155-164.
13. Bernardini, S. Exploring new directions for discovery learning. 165-182.
14. Kennedy, C. The CWIC project: Developing and using a corpus for intermediate Italian students. 183-192.
- Miceli, T.
15. Kübler, N. Linguistic concerns in teaching with language corpora: Learner corpora. 193-202.
16. Berglund, Y. The influence of external factors on learner performance. 205-215.
- Mason, O.
17. Leńko-Szymańska, A. How to trace the growth in learners' active vocabulary: A corpus-based study. 217-230.
18. Flowerdew, J. Computer-assisted analysis of language learner diaries: A qualitative application of word frequency and concordancing software. 231-243.
19. Lee, D. Genres, registers, text types, domains and styles: Clarifying the concepts and navigating a path through the BNC jungle. 247-292.
20. Gavioli, L. Some thoughts on the problem of representing ESP through small corpora. 293-303.
21. Thompson, P. Modal verbs in academic writing. 305-325.
22. Zanettin, F. CEXI: Designing an English-Italian translational corpus. 329-343.
23. Serpollet, N. Mandative constructions in English and their equivalents in French: Applying a bilingual approach to the theory and practice of translation. 345-359.
24. Claridge, C. Translating phrasal verbs. 361-373.

**TaLC 3.** Oxford, UK: Keble College. 24<sup>th</sup> – 27<sup>th</sup> July 1998.

Lou Burnard & Tony McEnery (eds). 2000. *Rethinking Language Pedagogy from a Corpus Perspective*. Frankfurt: Peter Lang.

1. Aston, G. Corpora and language teaching. 7-17.
2. Clear, J. Do you believe in grammar? 19-30.
3. Hoey, M. The hidden lexical clues of textual organisation: A preliminary investigation into an unusual text from a corpus perspective. 31-41.
4. Simpson, R. Methodological challenges of planning a spoken corpus with pedagogical outcomes. 43-49.
- Lucka, B.
- Ovens, J.
5. Collins, H. Materials design and language corpora: A report in the context of distance education. 51-63.

6. Foucou, P-Y. Kübler, N. A web-based environment for teaching technical English. 65-73.
7. Tribble, C. Genres, keywords, teaching: Towards a pedagogic account of the language of project proposals. 75-90.
8. Scott, M. Focusing on the text and its key words. 103-121.
9. Tono, Y. A computer learner corpus based analysis of the acquisition order of English grammatical morphemes. 123-132.
10. Cheng, W. Warren, M. The Hong Kong Corpus of Spoken English: Language learning through language description. 133-144.
11. Flowerdew, L. Investigating referential and pragmatic errors in a learner corpus. 145-154.
12. Eppler, E. The Interculture Project corpus: Data classification, access and the development of intercultural competence. 155-164.
13. Osborne, J. What can students learn from a corpus? Building bridges between data and explanation. 165-172.
14. Davies, M. Using multi-million word corpora of historical and dialectal Spanish texts to teach advanced courses in Spanish linguistics. 173-185.
15. Szakos, J. Producing and using corpora in Chinese language education. 187-192.
16. Hahn, A. Grammar at its best: The development of a rule- and corpus-based grammar of English tenses. 193-205.
17. Seidlhofer, B. Operationalizing intertextuality: Using learner corpora for learning. 207-223.
18. Bernardini, S. Systematising serendipity: Proposals for concordancing large corpora with language learners. 225-234.
19. Pearson, J. Surfing the internet: Teaching students to choose their texts wisely. 235-239.

**TaLC 2.** Lancaster, UK: Lancaster University. 9<sup>th</sup> – 12<sup>th</sup> August 1996.

Simon Botley, Julia Glass, Anthony McEnery & Andrew Wilson (eds). 1996. *Proceedings of TaLC 1996. UCREL Technical Papers 9*. Lancaster: University Centre for Computer Corpus Research on Language.

1. Kettemann, B. Concordancing in English language teaching. 4-16.
2. Magee, S. Rundell, M. The role of the corpus-based 'phrasicon' in English language teaching. 17-28.
3. Coniam, D. Using corpus word frequency data in the automatic generation of English language cloze tests. 29-44.
4. Barlow, M. Parallel texts in language teaching. 45-56.
5. Danielsson, P. Ridings, D. Corpus and terminology: Software for the translation program at Göteborgs Universitet, or getting students to do the work. 57-67.
6. Peters, C. Picchi, E. Biagini, L. Parallel and comparable bilingual corpora in language teaching and learning. 68-82.
7. Fernandez-Villanueva, M. Research into the functions of German modal particles in a corpus. 83-93.
8. Schmied, J. Encouraging students to explore language and culture in early modern English tracts. 94-107.
9. Gledhill, C. Science as a collocation: Phraseology in cancer research articles. 108-126.
10. Carne, C. Corpora, genre analysis and dissertation writing: An evaluation of the potential of corpus-based techniques in the study of academic writing. 127-137.
11. Jappy, T. Investigating grounding across narrative and oral discourse. 138-149.
12. Facchinetti, R. The exploration of diachronic English software by foreign language students. 150-159.
13. Bowden, P. Edwards, M. Knowledge extraction from corpora for pedagogical applications. 160-170.
14. Meyer, R. Okurowski, M. Hand, T. Using authentic corpora and language tools for adult-centred learning. 171-177.

15. Aston, G. The British National Corpus as a language learner resource. 178-191.
16. Burnard, L. Introducing SARA: An SGML-aware retrieval application for the British National Corpus. 192-202.
17. Pearson, J. Teaching terminology using electronic resources. 203-216.
18. Prince, V. A textual clues approach for generating metaphors as explanations by an intelligent tutoring system. 217-232.
- Ferrari, S.
19. Milton, J. Exploiting L1 and L2 corpora for computer assisted language learning design: The role of an interactive hypertext grammar. 233-243.
20. Dossena, M. Evaluating corpora: Are we asking the right questions? 244-253.
21. Hatzidakis, O. Corpus linguistics as an academic subject. 254-265.
22. Morley, J. A corpus-based description of headline grammar: The verb in headlines. 266-280.

**TaLC 1.** Lancaster, UK: Lancaster University. 10<sup>th</sup> – 13<sup>th</sup> April 1994.

a) Anne Wichmann, Steven Fligelstone, Tony McEnery & Gerry Knowles (eds). 1997.

*Teaching and Language Corpora*. Harlow: Addison Wesley Longman.

1. Leech, G. Teaching and language corpora: A convergence. 1-23.
2. Sinclair, J. Corpus evidence in language description. 27-39.
3. Mindt, D. Corpora and the teaching of English in Germany. 40-50.
4. Aston, G. Enriching the learning environment: Corpora in ELT. 51-64.
5. Minugh, D. All the language that's fit to print: Using British and American newspaper CD-ROMs as corpora. 67-82.
6. Gavioli, L. Exploring texts through the concordancer: Guiding the learner. 83-99.
7. Johns, T. Contexts: The background, development and trialling of a concordance-based CALL program. 100-115.
8. Wilson, E. The automatic generation of CALL exercises from general corpora. 116-130.
9. Dodd, B. Exploiting a corpus of written German for advanced language learning. 131-145.
10. Jones, R. Creating and using a corpus of spoken German. 146-156.
11. Ahmad, K. The role of corpora in studying and promoting Welsh. 157-172.
- Davies, A.
12. Peters, P. Micro- and macrolinguistics for natural language processing. 175-185.
13. Kettemann, B. Using a corpus to evaluate theories of child language acquisition. 186-194.
14. Knowles, G. Using corpora for the diachronic study of English. 195-210.
15. Wichmann, A. The use of annotated speech corpora in the teaching of prosody. 211-223.
16. Jackson, H. Corpus and concordance: Finding out about style. 224-239.
17. Louw, B. The role of corpora in critical literary appreciation. 240-251.
18. Renouf, A. Teaching corpus linguistics to teachers of English. 255-266.
19. Inkster, G. First catch your corpus: Building a French undergraduate corpus from readily available textual resources. 267-276.
20. King, P. Creating and processing corpora in Greek and Cyrillic alphabets on the personal computer. 277-291.
21. Hughes, G. Developing a computing infrastructure for corpus-based teaching. 292-307.

b) Andrew Wilson & Tony McEnery (eds). 1994. *Corpora in Language Education and Research*. UCREL Technical Papers 4. Lancaster: University Centre for Computer Corpus Research on Language.

1. Burstein, J. Parsing sentence fragments in computer-assisted test scoring. 1-8.
- Kaplan, R.
2. Esser, J. On the theoretical status of intonation in spoken corpora. 9-16.
3. van Halteren, H. Syntactic databases in the classroom. 17-28.
4. Kirk, J. Teaching and language corpora: The queen's approach. 29-51.
5. Kita, K. Automatically extracting collocations from corpora. 53-64.
- Kato, Y.

- Omoto, T.  
Yano, Y.
6. Milton, J. A corpus-based online grammar and writing tool for EFL learners: A report on work in progress. 65-77.
7. Paulussen, H. Automating vocabulary teaching material: A step-by-step experience. 79-85.  
Deville, G.
8. Quinn, D. Linguistic modelling for a corpus-based CALL system. 87-98.  
Quinn, A.
9. Zanettin, F. Parallel words: Designing a bilingual database for translation activities. 99-111.