# Integrating corpus tools on intensive CELTA courses

*Ben Naismith*

*In the fields of Corpus Linguistics and Teacher Education, there is a substantial body of research relating to corpus applications in the classroom. However, the majority of such work relates to teacher training in tertiary education. In this article, a research project is described in which corpus training was included as part of a Certificate in English Language Teaching to Adults (CELTA), the popular four-week teacher training programme. Specifically, the project focuses on web-based corpus tools that provide frequency information. Observations and evidence are offered from both the trainees' and trainers' perspectives, and the subsequent analysis reveals a number of trends. Overall, it is argued that corpus tools are a valuable training resource on CELTA programmes in terms of developing language awareness and assisting in lesson planning, but that expecting trainees to use corpus tools during teaching practice remains unrealistic. In the ensuing discussion, practical considerations are provided for the integration of corpus tools into CELTA programmes.*

## Introduction

In ELT, a premium is placed on the ability of the teacher to act as a model of the language. From assessment criteria on Cambridge teacher training awards to discriminatory 'native speakerism'[1] in job advertisements, the idea that the teacher should be the exemplar for their learners is prevalent. And yet, although teachers-in-training (trainees) taking courses like the Certificate in English Language Teaching to Adults (CELTA) are typically expert speakers,[2] often expert speakers only have partial knowledge of a language (Krishnamurthy 2001: 172).

From personal experience, it seems that when entering teacher training programmes, many trainees are sensitive to the fact that their declarative knowledge of grammar is limited. In contrast, there is often a 'blind spot' for their incognizance of lexis and that their lexical awareness requires development too. Consequently, trainees regularly dedicate a disproportionate amount of their time to studying grammar while neglecting to consider key elements of lexis, especially collocation. This issue is only further compounded by the popular ELT coursebooks which are commonly used as resources on CELTA and predominantly adhere to a grammar syllabus.

In order to address this discrepancy, one option is to incorporate corpora, i.e. 'large collections of naturally occurring discourse' (Chambers 2010: 345), into teacher training. To date, corpora remain underused in much

teacher education but are arguably an untapped resource with true pedagogical potential.

**Corpora in ELT**
Corpora and language learners

Conventionally, corpora use in ELT is divided into two categories, 'direct use' and 'indirect use' (Römer 2006: 124), although numerous naming conventions exist. No matter the designation, the defining features remain constant: for direct use, it is that corpus data be utilized by students as a resource in a hands-on manner. In contrast, in indirect use, corpus data are accessed by experts, for example writers and teachers, to create coursebooks, dictionaries, and other materials.

At the initial vanguard of direct corpus use was Tim Johns, the originator of the inductive approach he named data-driven learning (DDL). Emerging in the mid-1980s, Johns' influential work brought DDL to the forefront of corpus use in the classroom (Boulton 2010). Of key importance, Johns objected to corpus input becoming 'virtually invisible to the learner' (ibid.: 1) and pushed learners towards discovery of language for themselves, famously describing DDL as an '[attempt] to cut out the middleman' (Johns 1991: 30). In essence, DDL advocates that learners access corpus data in the classroom, typically through computer software, to notice features of language and to draw conclusions about its usage.

Perhaps the most commonly cited objection to DDL is the authenticity of concordances, i.e. that concordance lines are isolated, partial sentences lacking context (cf. Widdowson 2000). However, as Chambers (op.cit.) has noted, corpus software does typically contain links to the complete source texts. Other issues raised relate to proficiency with technology and corpus software, for example that corpora can be too technically challenging for teachers and language learners, although training and support options have been proposed (cf. Cheng, Warren, and Xu 2003). In addition to the above criticisms, a survey of DDL research reveals a paucity of quantitative data and large-scale research with many studies, including the one described in this article, based on small samples and self-reporting.

Corpora and trainees

Within teacher education, the findings in relation to direct corpus use and language learners are equally applicable to novice teachers. That DDL might be useful to trainees is unsurprising given that the goals of both language learners and trainees are identical in many regards; for example, both groups are seeking to develop their own language awareness. Supporting this premise, studies have shown possible uses of DDL with trainees in terms of improving language awareness, skills development, and critical thinking (cf. Coniam 1997). Nonetheless, in terms of corpus use in teacher education, the current literature presents a clear consensus best summarized by Römer (op.cit.: 121), who writes that despite the potential benefits, 'corpora and corpus methods play a very minor role in EFL initial teacher training at universities'. To remedy this shortcoming, many researchers push to see further corpora training at all levels of teacher education, whether as part of university programmes, teacher-training courses, or in-service professional development (cf. Frankenberg-Garcia 2012: 35).

*Ben Naismith*

Surveying both direct and indirect corpus uses with learners and teachers, we see that the boundaries between the two are often blurred; as Boulton (op.cit.: 28) notes, DDL is not an 'all-or-nothing process'. Tasks falling into this grey area include using teacher-edited corpus texts, providing students with corpus search terms, or displaying corpus information to aid with clarification. When considering teacher training, it is challenging then to distinguish whether DDL is being strictly adhered to as trainees simultaneously take on the roles of learner and teacher as part of their studies. In 'softening' the total freedom of strong-form DDL, many researchers see a justifiable compromise to overcome corpora resistance (Chambers op.cit.), a position I adopt in my research.

Corpus web tools

In many ELT contexts, the widespread presence of web tools is undeniable. As such, it is imperative to consider how these technologies can be incorporated into teacher-training programmes. At present, for virtually all corpus software, the most commonly used tool is the concordancer, which provides a list of authentic examples of words, i.e. concordances (Krishnamurthy op.cit.: 178). Allowing close examination of how words are used, concordances are an excellent insight into naturally occurring examples of language; however, use of concordancers also typically requires some expertise and training, and the results may be challenging to interpret. In contrast, frequency information is simple to decode and its utility is immediately apparent to trainees, who are able to apply the rudimentary, if imperfect, equation, 'most frequent = most important to learn' (Leech 1997: 16). Below, examples of two different types of corpus tools are highlighted which can display such information: frequency trackers and simplified corpus interfaces.

**Frequency trackers**
One simple yet powerful frequency tracker is Google Ngram Viewer, which uses a clear line graph to display word usage over time, based on the massive Google Books corpus. Users need only input words or phrases and Ngram Viewer instantly creates an interactive graph (Figure 1). In addition to its basic functions, more advanced searches may include parts of speech, genres, and other metadata. Despite the ease of access and effective graphics, the ELT community has yet to research Ngram Viewer's
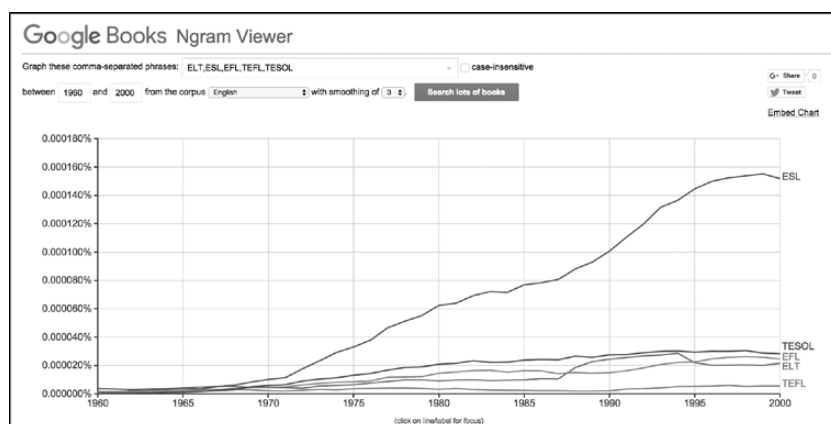
FIGURE 1
Ngram Viewer: common language teaching acronyms

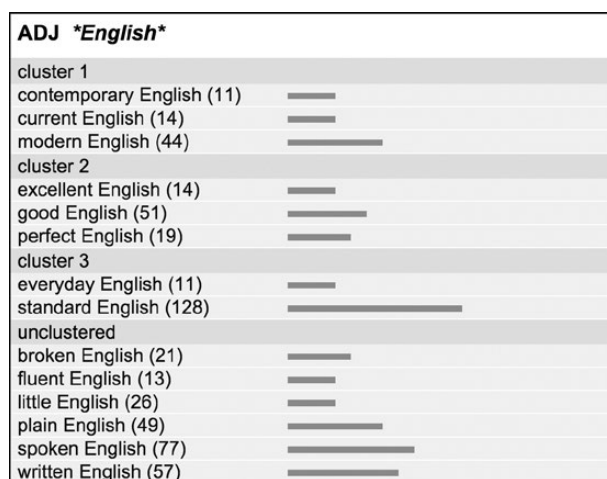| ADJ *English* | |
|---|---|
| **cluster 1** | |
| contemporary English (11) | ▬ |
| current English (14) | ▬ |
| modern English (44) | ▬▬ |
| **cluster 2** | |
| excellent English (14) | ▬ |
| good English (51) | ▬ |
| perfect English (19) | ▬ |
| **cluster 3** | |
| everyday English (11) | ▬ |
| standard English (128) | ▬▬▬ |
| **unclustered** | |
| broken English (21) | ▬ |
| fluent English (13) | ▬ |
| little English (26) | ▬ |
| plain English (49) | ▬▬ |
| spoken English (77) | ▬▬ |
| written English (57) | ▬▬ |

FIGURE 2
Just-the-Word: adjective collocations with 'English'

potential in any great depth; at the time of writing there are no known publications relating Ngram Viewer to ELT.

## The study

### Simplified corpus interfaces

In contrast to Ngram Viewer's pure focus on frequency, sites like www.just-the-word.com offer frequency information in addition to concordancing options. However, what sets Just-the-Word apart from its concordancing counterparts is its interface; intuitive and accessible, this 'graded' corpus tool can be readily used by novices. Based on the British National Corpus, Just-the-Word merely requires users to enter a search term. Upon doing so, the site distils the information down to its essentials in terms of frequency, i.e. a list of the most common collocates, sorted by pattern. Instead of bombarding users with statistics, Just-the-Word displays bars to indicate relative frequency: the longer the bar, the more frequent the collocation. Thus, even at a glance, a learner or teacher can instantly determine the most common collocates (Figure 2).

### Research context

The project reported in this article was a two-month long investigation that took place at IH Vancouver, a private language institute in Western Canada. The project sought to examine more closely how the benefits from corpora could be introduced at a critical time in a teacher's formation, during their initial pre-service training. To help in addressing gaps in the existing literature, the study focused on those aspects of classroom corpus use which have to date been under-researched, that is non-concordancing corpus tools and short pre-service training courses. Specifically, the research attempted to answer the following two questions:

1  On initial teacher-training courses such as CELTA, do trainees perceive the inclusion of frequency-focused corpus web tools to be (a) of interest, and (b) beneficial, in developing their own language awareness and teaching?
2  Having received training on the use of these corpus tools, and given the opportunity, will trainees use them in their own lesson planning and teaching practice?

At IH Vancouver, eight CELTA courses typically run every year. These courses have an intensive four-week format and contain a practicum

component of eight assessed lessons. In the study, which covered two such courses, there were 16 trainees in total, with diverse cultural and linguistic backgrounds. Consent was obtained from all participants prior to the courses and all necessary steps were taken to ensure the ethicality of the research; participation was completely voluntary and did not affect trainees' assessment in any way.

Methodology

In selecting the methodological framework, Action Research (AR) was deemed the most suitable as it is commonly employed to increase understanding of, and bring about change to, classroom practices, including teacher training (Richards 1996: 12). Furthermore, the narrow scope of the project lent itself to AR, which 'typically involves small-scale investigative projects in the teacher's own classroom' (ibid.).

Although many conceptions of AR exist, for this research the simple cyclical model (Figure 3) was selected, thereby allowing for a sequence in which the findings could be applied to future cycles of investigations. Within this model (Hudson, Owen, and van Veen 2003), there are four broad stages in which the researcher

- plans a course of action;
- acts to implement the plan;
- observes the effects of the action; and finally
- reflects on the previous stages, with a view to possibly repeating the cycle.

For this research, the planning stage started with the observation that corpora were not being used on CELTA programmes, which led naturally to the project's research questions and literature review. The acting stage then entailed providing the corpus training and the data collection. Intertwined with the acting stage was the observing stage; having introduced the corpus tools, tutors observed participants' teaching practice and made field notes. Finally, post-course, the data analysis and writing of this article constituted the reflecting stage, with opportunities for the entire cycle to be refined and repeated on future courses.

To provide training related to corpus tools, one 75-minute input session was delivered midway through the course. This session, focusing on lexis, aimed to deepen language awareness of multi-word units and to provide ideas and tools relating to corpora. In addition, a written language-analysis assignment was given in the first week as a diagnostic test to see how trainees would use corpus tools prior to receiving corpus training. As part of this assignment, trainees were required to find common collocates of specific lexical items and were provided with a short list of possible online resources, including the web tools described earlier in this article.
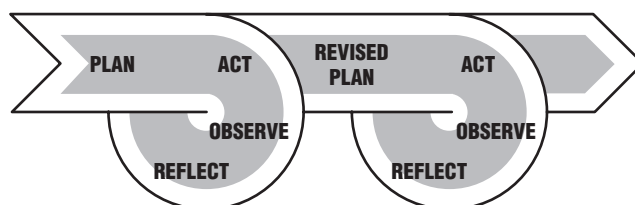


FIGURE 3
Action Research cycle
(Hudson *et al.* op.cit.)

| Data collection and analysis | A combination of quantitative and qualitative data was compiled throughout the research using three tools: a questionnaire, field notes, and an observation table. By using a variety of collection methods, the aim was to gather data from the perspective of the participants as well as the observers and in relation to multiple aspects of the course. The questionnaire, given on the final day, was the primary measurement tool for collecting participants' self-reported experiences of using and learning about corpus tools. To complement this self-reported data, field notes and an observation table were completed by the course tutors based on the trainees' teaching practice. |

A multifaceted analysis of the data was then conducted. The goal in this case was not for one type of data to dominate, but for the quantitative data to reveal general patterns, with the qualitative data then adding contextual information from the viewpoint of the participants (Bryman and Burgess 1994: 222). In analysing the quantitative data, the raw figures were used to uncover statistical trends and were supplemented with chi-square tests to determine correlation between the variables under investigation. In contrast, analysis of the qualitative data required greater interpretation to extract common patterns, with trainees' comments and tutors' fieldnotes helping to support the statistical evidence.

**Findings**

Having analysed the collected results, five main trends emerged, providing insight into direct corpus use on the CELTA.

Trend 1: high interest and perceived benefits

Interest in the corpus tools was consistently high, as evidenced by the questionnaire answers and tutor observations, and regardless of whether trainees even used the corpus tools. Such a result is consistent with previous findings in other studies which also relied on self-reporting (cf. Cheng *et al.* op.cit.) and does not appear to have been affected by the comparatively shorter length of the course.

Typical reasons given for trainees' interest in corpus tools were invariably linked to the perceived benefits of using them during lessons, for example 'It's a great way to add something new to the class', and for becoming more knowledgeable about authentic language production, for example 'These seem like useful tools for understanding language use and providing clear quantifiable evidence'. Similarly, there was a high degree of interest in future use with most trainees claiming that they were likely to incorporate corpus tools in their upcoming teaching careers. Reasons given for these predictions again related to their utility as a teaching tool, alternately described as 'helpful', 'awesome', and 'a great resource'.

Trend 2: low familiarity

Superficially, low familiarity with corpus tools may seem inevitable considering that the CELTA is an initial teaching qualification. However, CELTA trainees do often have pre-existing teaching experience, and in the case of this project, 50 per cent had worked in ELT prior to taking the course. Yet, despite this level of experience, familiarity with corpus tools was consistently low, with only one trainee having actually used them. These results are likely indicative that both the general public and members of the ELT industry tend to lack exposure to corpus tools. What is not known, but would be of interest, is the percentage of CELTA trainers who are equally unfamiliar with corpora.

*Ben Naismith*

| Trend 3: low usage | The data showed uniformly low usage of corpus tools for trainees on both courses. As these tools were completely optional and did not affect their grades, there were no extrinsic motivating factors to use the tools in their planning or teaching. In total, approximately half of candidates (56 per cent) used the tools in 9 per cent of the lessons. However, a further breakdown of this figure indicates that prior to the corpus training midway through the course, there was only one isolated occurrence of a trainee using a corpus tool (2 per cent of lessons), whereas after the training this figure jumped to 17 per cent. As such, it seems that corpus training, combined with a greater level of experience in the classroom, is essential if corpus tools are to be used voluntarily. Statistically, when using a chi-square test, the $p$ value of this correlation between the second half of the course and corpus tools being used was significant at $p < 0.01$. In contrast, other possible variables which were analysed failed to explain who would use the corpus tools, yielding insignificant results. These variables included trainees' final course grades, previous exposure to corpus tools, and reported levels of interest. |
|---|---|
| Trend 4: language analysis as primary use | Among lessons that included corpus tools, there was a clear disparity between usage in lesson planning as compared to teaching practice, with nearly all uses of corpus tools occurring before (92 per cent), rather than during (8 per cent), the lesson. Trainees indicated that this preference stemmed from insecurities about using the technology while teaching, electing instead to consult it at home in an untimed, unassessed setting either 'to check collocations and chunks', '[to decide] which lexis to focus on', or '[to help] clarify best language use'. Other trainees pointed to the intensive context, with one feeling 'too stressed out to be creative'. Such a trend indicates that trainees found DDL to be a useful approach when considering themselves as learners, but in their roles as teachers they did not extend DDL opportunities to their own students, opting for them to only receive the benefits of indirect corpus use. |

Even within planning, one usage of corpus tools, language analysis, was noticeably more frequent than others. In the majority of cases, language analysis, especially lexical analysis, was the motivating factor for consulting corpora, accounting for 67 per cent of uses when planning, with grammar analysis and materials design a distant second at 25 per cent each.[3] Specific examples of such uses included comparing the frequency of 'started work as' versus 'started working as' or determining whether 'traveled' or 'travelled' was the more common spelling. In terms of justification, trainees posited that corpus tools could most easily be used as a lexical reference rather than as the basis for materials, especially considering the wealth of existing activities available online and in the school's library. In contrast, trainees found few readily available sources which provided information about lexis relating to collocations or frequency.

| Trend 5: choice of corpus tools | In selecting corpus tools, trainees were consistent in their preferences. At first, when given the language-focused assignment, trainees largely opted for Just-the-Word. In contrast, after receiving corpus training, Ngram Viewer became the primary choice. Although this change in preference may appear arbitrary, the most likely explanation relates to ease of use: |
|---|---|

whichever tool could more easily be used to achieve the desired result was the one selected. Thus, for the language assignment which required learners to find unknown collocations, Just-the-Word was simpler as trainees needed to only type in the word, click, and read the results. In contrast, when comparing two known terms, for example 'cellphone' and 'mobile phone', trainees found it simpler to use Ngram Viewer as no metalanguage was necessary and the resulting visual graph was both appealing and easy to interpret. As this latter use tended to be more common when planning lessons, Ngram Viewer was the most popular tool (83 per cent of occurrences), followed by Just-the-Word (42 per cent).

## Discussion
### The CELTA context

In discussing the above trends, certain key aspects of 'the CELTA context' are worth considering. First, in terms of length, the majority of CELTA courses are four weeks, and, as a result, intensive in nature due to the range of skills and knowledge that trainees are expected to demonstrate. Whether this expectation is reasonable has been debated, but the reality remains and is reflected in the 42 criteria which comprise the teaching practice assessment (Cambridge English 2015).

As a result of this context, a number of issues exist relating to standardizing the implementation of corpus tools on CELTA. For one, as CELTA is a truly global qualification, it is impossible to guarantee that in all regions of the world there would be the necessary access to technology; although open-access corpora like the Corpus of Contemporary American English (COCA) have removed one potential barrier, reliable internet and computer availability are not universal. Likewise, in certain training contexts, candidates may not possess the requisite level of computer literacy required to effectively use corpus tools. Consequently, if corpus use were to become a mandatory component of CELTA, it could disqualify potential centres and candidates. In terms of course tutors, foreseeable issues also arise, namely, that adding a standardized corpus element would prove logistically problematic as there is undoubtedly great variance amongst tutors' familiarity with corpora.

Ultimately, both of these hypothetical issues can be seen to relate to prescriptive use of corpora and corpus tools. In taking such a stance, we would run the risk of discouraging teachers from using corpus tools, thereby widening the gap between linguistic theory and classroom practice. Alternatively, '[instead] of asking "What can a teacher do with a corpus?" we might ask "What can a corpus do for a teacher?"' (Frankenberg-Garcia op.cit.: 36), taking only those steps which we feel would add value to trainees' CELTA experience without also imposing additional burden on them.

### Recommendations

In adhering to this guiding principle, I offer one realistic recommendation for CELTA at each of the local and global levels to try to help teachers along the road to expertise, where 'one must be more *knowledgeable*, be more *efficient*, and have better *insight* than non-experts' (O'Keeffe and Farr 2003: 392; italics in original).

*Ben Naismith*

### Local recommendation: encouraging unassessed usage

To further encourage use of corpus tools at IH Vancouver, the relevant input session will now take place in the computer lab rather than the classroom. As these facilities are available, it will allow for more hands-on use of corpus tools, in line with the suggestion that trainees 'should be encouraged to start with free online corpora, which they can try out with no strings attached' (Frankenberg-Garcia op.cit.). Likewise, as part of the pre-course task sent to trainees, an activity will be included to prompt an introductory exploration of basic corpus information. In doing so, familiarization of corpus tools may take place without impinging on the existing course timetable.

### General recommendation: amending assessment criteria

At a broader level, I recommend small changes to the CELTA syllabus in order to reflect an increased emphasis on evidence-based lexical awareness. Such changes would not require drastic amendments as courses like CELTA already assess language-knowledge criteria. Specifically, in relation to lexis, there are sub-criteria relating to meaning, form, and pronunciation (Cambridge English op.cit.). Nevertheless, these points deal almost exclusively with form at the word level rather than larger lexical patterns, with minimal mention of collocation. As such, minor alterations could reasonably be added to increase the emphasis on these key aspects of lexis.

In including 'frequency' as a key component of lexis, there would thereby be greater incentive for trainers and trainees to analyse lexis for this characteristic, one which is often overlooked by textbooks. And by extension, resources like corpus tools which provide such information could be employed to combat a reliance on the notoriously unreliable intuition of 'expert' users (Krishnamurthy op.cit.). Furthermore, a focus on frequency would allow for the scope of the language awareness training to include lexicogrammatical patterns rather than artificially dividing language into the mutually exclusive categories of 'vocabulary' and 'grammar'. Similarly, by supplementing 'context of situation' with 'appropriacy', there would be a more direct link between teacher/learner language (including the frequency of items) and the effect of lexical choice on communication (see amendments in bold in Figure 4).

## Conclusions

Returning to the original research questions, the findings do suggest certain conclusions, bearing in mind the small sample size and specific context of the action research. For the first question, it seems clear that this CELTA-based

FIGURE 4
Proposed amendments to CELTA syllabus (Cambridge English op.cit.)

| 2.3 | Lexis<br>Word formation, meaning and use in context | a. understand basic principles of word formation and lexical meaning, for example:<br>– meaning and definition<br>– pronunciation<br>– spelling<br>– affixation and compounding<br>– synonymy and hyponymy<br>– **frequency**<br><br>b. understand the effect on word choice of factors such as:<br>– co-text (e.g. collocation)<br>– context of situation (style, **appropriacy**) |

study matches other university-based studies in that trainees are interested in, and perceive, the benefits of corpus tools. Unlike other studies, however, there were no issues relating to corpus tool training or technology, a divergence it is possible to attribute to the restricted focus on easy-to-use, frequency-focused tools rather than more complex corpus interfaces or concordancers.

In response to the second research question relating to actual use of corpus tools, the answer is two-fold: overall, usage of corpus tools was not overwhelming, yet nevertheless significant given the context of this CELTA course, especially in the second half of the course after minimal corpus training. More notably, the evidence suggests that use of corpus tools in terms of developing trainees' language awareness and assisting in lesson planning is a realistic and powerful option. However, expecting trainees to then apply the techniques of DDL to their own assessed teaching practice is idealistic and unlikely to occur.

Following on from the current investigation, through future cycles of action research I hope to validate the results, with a view to maximizing the usefulness of corpora in my training context. Likewise, in order to deepen the understanding of this intersection of teacher education and corpus linguistics, I would encourage further related studies involving alternative corpus tools and corpora, or use of corpus tools on other teacher training programmes such as the Cambridge DELTA modules.

*Final version received August 2016*

## Notes
1 The belief that ideal teachers are native-speakers due to their spoken language proficiency and supposedly more prestigious language norms.
2 'Expert speakers' includes native and non-native speakers.
3 Some figures relating to usage total more than 100 per cent as options are not mutually exclusive.

## References
**Boulton, A.** 2010. 'Data-driven learning: on paper, in practice' in T. Harris and M. Jaén (eds.). *Corpus Linguistics in Language Teaching*. Frankfurt: Peter Lang.

**Bryman, A.** and **R. G. Burgess** (eds.). 1994. *Analyzing Qualitative Data*. Abingdon: Routledge.

**Cambridge English**. 2015. *CELTA (Certificate in Teaching English to Speakers of Other Languages): Syllabus and Assessment Guidelines* (fourth edition). Available at http://www.cambridgeenglish.org/images/21816-celta-syllbus.pdf (accessed on 20 January 2016).

**Chambers, A.** 2010. 'What is data-driven learning?' in A. O'Keeffe and M. McCarthy (eds.). *Routledge Handbook of Corpus Linguistics*. London: Routledge.

**Cheng, W.**, **M. Warren**, and **X. Xu.** 2003. 'The language learner as language researcher: putting corpus linguistics on the timetable'. *System* 31/2: 173–86.

**Coniam, D.** 1997. 'A practical introduction to corpora in a teacher training language awareness programme'. *Language Awareness* 6/4: 199–207.

**Frankenberg-Garcia, A.** 2012. 'Integrating corpora with everyday language teaching' in J. Thomas and A. Boulton (eds.). *Input, Process and Product: Developments in Teaching and Language Corpora*. Brno: Masaryk University Press.

**Hudson, B.**, **D. Owen**, and **K. van Veen.** 2003. 'Working on educational research methods with Masters students in an international online learning community'. Paper presented at European Conference on Educational Research, University of Hamburg, 17–20 September.

**Johns, T.** 1991. 'Should you be persuaded: two samples of data-driven learning materials'. *English Language Research Journal* 4/1: 1–16.

**Krishnamurthy, R.** 2001. 'Size matters: creating dictionaries from the world's largest corpus'. *KOTESOL Proceedings 2000*. Taegu: KOTESOL. Available at https://koreatesol.org/content/kotesol-proceedings-2000.

**Leech, G.** 1997. 'Teaching and language corpora: a convergence' in A. Wichmann, S. Fligelstone, A. McEnery, and G. Knowles (eds.). *Teaching and Language Corpora*. London: Longman.

**O'Keeffe, A.** and **F. Farr.** 2003. 'Using language corpora in initial teacher education: pedagogic issues and practical applications'. *TESOL Quarterly* 37/3: 389–418.

**Richards, J. C.** 1996. *Reflective Teaching in Second Language Classrooms*. Cambridge: Cambridge University Press.

**Römer, U.** 2006. 'Pedagogical applications of corpora: some reflections on the current scope and a wish list for future developments'. *Zeitschrift für Anglistik und Amerikanistik* 54/2: 121–34.

**Widdowson, H.** 2000. 'The limitations of linguistics applied'. *Applied Linguistics* 21/1: 3–25.

## The author

**Ben Naismith** is a teacher, teacher trainer, and aspiring researcher. Since 2002 he has worked in a range of contexts, most notably in Costa Rica, Thailand, the UAE, and Canada. He holds an MSc in Applied Linguistics from Aston University and a DELTA from Cambridge English. His current research interests include teacher education, corpus linguistics, approaches and methods in language learning, and teacher language awareness. Twitter: @BenNaismithELT.
**Email:** bennaismithelt@gmail.com