

Metadata of the chapter that will be visualized in SpringerLink

Book Title	Open Educational Resources in Higher Education	
Series Title		
Chapter Title	Reflections on Remixing Open Access Content into Open Educational Resources: A New Paradigm for Sustainable Data-Driven Language Learning Systems Design in Higher Education	
Copyright Year	2023	
Copyright HolderName	The Author(s), under exclusive license to Springer Nature Singapore Pte Ltd.	
Corresponding Author	Family Name	Fitzgerald
	Particle	
	Given Name	Alannah
	Prefix	
	Suffix	
	Role	
	Division	Department of Computer Science
	Organization	The University of Waikato
	Address	Te Whare Wananga o Waikato, Gate 1, Knighton Road, Hamilton, 3240, New Zealand
	Division	
	Organization	Durham University
	Address	Durham, UK
	Email	alannah.fitzgerald@waikato.ac.nz
	ORCID	http://orcid.org/0000-0003-0392-2740
Author	Family Name	Wu
	Particle	
	Given Name	Shaoqun
	Prefix	
	Suffix	
	Role	
	Division	Department of Computer Science
	Organization	The University of Waikato
	Address	Te Whare Wananga o Waikato, Gate 1, Knighton Road, Hamilton, 3240, New Zealand
	Email	shaoqun.wu@waikato.ac.nz
	ORCID	http://orcid.org/0000-0001-9566-005X
Author	Family Name	König
	Particle	
	Given Name	Jemma
	Prefix	
	Suffix	
	Role	
	Division	Department of Computer Science
	Organization	The University of Waikato

Address Te Whare Wananga o Waikato, Gate 1, Knighton Road, Hamilton, 3240, New Zealand
Email jemma.konig@waikato.ac.nz

Author

Family Name **Shaw**
Particle
Given Name **Steven**
Prefix
Suffix
Role
Division Department of Education
Organization Concordia University
Address Faubourg Ste-Catherine Building, 1610 St. Catherine West, Montreal, QC, H3G 1M8, Canada
Email steven.shaw@concordia.ca
ORCID <http://orcid.org/0000-0002-0719-6670>

Author

Family Name **Witten**
Particle
Given Name **Ian H.**
Prefix
Suffix
Role
Division Department of Computer Science
Organization The University of Waikato
Address Te Whare Wananga o Waikato, Gate 1, Knighton Road, Hamilton, 3240, New Zealand
Email ihw@waikato.ac.nz
ORCID <http://orcid.org/0000-0001-6428-8988>

Abstract

This chapter presents a new paradigm for sustainable data-driven language learning systems design in higher education that draws on qualitative reflections spanning a decade (2012–2022) with stakeholders from an ongoing global research study with the FLAX (Flexible Language Acquisition) and F-Lingo projects at the University of Waikato in Aotearoa New Zealand (Fitzgerald (2019) A new paradigm for open data-driven language learning systems design in higher education; König et al. (2022) *Smart CALL*). Design considerations are presented for remixing domain-specific open access content into Open Educational Resources (OER) for academic English language provision across formal and non-formal higher education contexts. Primary stakeholders in the research collaboration include the following three groups: (1) Knowledge organisations that provide open access to academic content—libraries and archives, including the British Library and the Oxford Text Archive, universities in collaboration with MOOC providers and the CORE (COnnecting REpositories) open access aggregation service at the UK Open University; (2) Researchers who mine and remix academic content into corpora and open data-driven language learning systems—converging from the fields of open education, computer science and applied corpus linguistics; (3) Knowledge users who re-use and remix academic content into OER—English for Academic Purposes (EAP) practitioners from university language centres. Automated content analysis was carried out on a corpus of interview and focus discussion data with the three stakeholder groups in this research. We discuss themes arising from the research data that reflect the different stakeholders’ experiences of remixing open access research content that has been produced within the academy for re-use as open educational content for teaching and learning features of academic language within open data-driven language learning systems. These open learning systems have been specifically designed to scale with OER expansion and traction in mind for their sustainable uptake both within and beyond the brick and mortar of the traditional university. The new paradigm presented in this chapter challenges, as the OER movement must, established business models and deeply embedded cultural or institutional norms that present obstacles to OER expansion and traction and the sustainability of the movement. One persistent challenge concerns the lack of open education policy across the higher education sector for full open

access (for use, modification, adaptation) via Creative Commons licensing to content produced within the academy. Thus, while this research has theoretical and practical implications in applied linguistics, computer science, language teaching and learning and open education, more generally, it also has significant cultural, business model and policy implications for higher education.

Keywords
(separated by '-')

Data-driven learning - Design-based research - English for academic purposes (EAP) - Higher education - Massive open online courses (MOOCs) - Open access - Open educational practices - Open educational resources (OER) - Systems design

Chapter 6

Reflections on Remixing Open Access Content into Open Educational Resources: A New Paradigm for Sustainable Data-Driven Language Learning Systems Design in Higher Education



Alannah Fitzgerald , Shaoqun Wu , Jemma König, Steven Shaw , and Ian H. Witten 

1 **Abstract** This chapter presents a new paradigm for sustainable data-driven language
 2 learning systems design in higher education that draws on qualitative reflections
 3 spanning a decade (2012–2022) with stakeholders from an ongoing global research
 4 study with the FLAX (Flexible Language Acquisition) and F-Lingo projects at the
 5 University of Waikato in Aotearoa New Zealand (Fitzgerald (2019) A new paradigm
 6 for open data-driven language learning systems design in higher education; König
 7 et al. (2022) *Smart CALL*). Design considerations are presented for remixing domain-
 8 specific open access content into Open Educational Resources (OER) for academic
 9 English language provision across formal and non-formal higher education contexts.
 10 Primary stakeholders in the research collaboration include the following three groups:
 11 (1) Knowledge organisations that provide open access to academic content—libraries

A. Fitzgerald (✉) · S. Wu · J. König · I. H. Witten
 Department of Computer Science, The University of Waikato, Te Whare Wananga o Waikato,
 Gate 1, Knighton Road, Hamilton 3240, New Zealand
 e-mail: alannah.fitzgerald@waikato.ac.nz

S. Wu
 e-mail: shaoqun.wu@waikato.ac.nz

J. König
 e-mail: jemma.konig@waikato.ac.nz

I. H. Witten
 e-mail: ihw@waikato.ac.nz

A. Fitzgerald
 Durham University, Durham, UK

S. Shaw
 Department of Education, Concordia University, Faubourg Ste-Catherine Building, 1610 St.
 Catherine West, Montreal, QC H3G 1M8, Canada
 e-mail: steven.shaw@concordia.ca

© The Author(s), under exclusive license to Springer Nature Singapore Pte Ltd. 2023
 J. Olivier and A. Rambow (eds.), *Open Educational Resources in Higher Education*,
 Future Education and Learning Spaces, https://doi.org/10.1007/978-981-19-8590-4_6

1

12 and archives, including the British Library and the Oxford Text Archive, universi-
13 ties in collaboration with MOOC providers and the CORE (COnnecting REposito-
14 ries) open access aggregation service at the UK Open University; (2) Researchers
15 who mine and remix academic content into corpora and open data-driven language
16 learning systems—converging from the fields of open education, computer science
17 and applied corpus linguistics; (3) Knowledge users who re-use and remix academic
18 content into OER—English for Academic Purposes (EAP) practitioners from univer-
19 sity language centres. Automated content analysis was carried out on a corpus of inter-
20 view and focus discussion data with the three stakeholder groups in this research. We
21 discuss themes arising from the research data that reflect the different stakeholders’
22 experiences of remixing open access research content that has been produced within
23 the academy for re-use as open educational content for teaching and learning features
24 of academic language within open data-driven language learning systems. These
25 open learning systems have been specifically designed to scale with OER expansion
26 and traction in mind for their sustainable uptake both within and beyond the
27 brick and mortar of the traditional university. The new paradigm presented in this
28 chapter challenges, as the OER movement must, established business models and
29 deeply embedded cultural or institutional norms that present obstacles to OER expansion
30 and traction and the sustainability of the movement. One persistent challenge
31 concerns the lack of open education policy across the higher education sector for full
32 open access (for use, modification, adaptation) via Creative Commons licensing to
33 content produced within the academy. Thus, while this research has theoretical and
34 practical implications in applied linguistics, computer science, language teaching
35 and learning and open education, more generally, it also has significant cultural,
36 business model and policy implications for higher education.

37 **Keywords** Data-driven learning · Design-based research · English for academic
38 purposes (EAP) · Higher education · Massive open online courses (MOOCs) ·
39 Open access · Open educational practices · Open educational resources (OER) ·
40 Systems design

41 6.1 Introduction

42 In this chapter section, we present a new research paradigm for sustainable data-
43 driven language learning systems design in higher education. This research paradigm
44 provides a theoretical and conceptual framework supported by a review of the relevant
45 literature from intersecting fields in this research. As we progress through the chapter,
46 the paradigm will be unpacked in greater detail in the subsequent sections as we drill
47 down into the specifics of the research contexts, materials and methods that have been
48 employed in the research with the three stakeholder groups. Reflections on remixing
49 open access content into OER for teaching and learning features of academic English,
50 along with the affordances and challenges encountered by the stakeholder groups,
51 will be presented in the final discussion section of this chapter.

6.1.1 Research Paradigm

A basic premise underpinning the new research paradigm presented in this chapter is that open data-driven language learning systems design as an approach is learner-centric and operates with the interface to the learner. Whether the learner is operating fully online in non-formal or informal learning mode or in a blended modality that is based both within and beyond the formal university language classroom, this approach requires that the tools and interfaces, and indeed the academic language corpora, be openly accessible and remixable for development or adaptation to meet this specific learner requirement. This method is different from existing Data-Driven Learning (DDL) approaches which assume specialised knowledge or experience with DDL tools, interfaces and strategies, operating on mostly inaccessible corpora in terms of cost or design, or assuming training to, hopefully, compensate for this lack of knowledge and experience (Fitzgerald, 2019; Pérez-Paredes et al., 2018).

The term DDL was coined by corpus linguistics and EAP pioneer, Tim Johns, to refer to a computer-driven language learning support approach with linguistic content that has been automatically analysed, enriched and transformed into a data-mined resource that learners can browse and query (Johns, 1991a). Johns envisioned every language learner as “a Sherlock Holmes” with direct access to the evidence of real-world language data (Johns, 2002, p. 108). In a similar vein to contemporary advocates for using and developing a broad spectrum of data literacies with open data in higher education (Atenas et al., 2015), Johns also envisioned DDL as developing data literacies for understanding and interpreting linguistic data for direct applications in language learning, specifically in the context of higher education (Johns, 2002; Pérez-Paredes et al., 2018).

From a research and development (R&D) standpoint, the paradigm presented here also operates with the interface to knowledge organisations (universities, libraries, archives) and researchers who are engaging with open educational practices to push at the parameters of open policy for the non-commercial re-use and remix of authentic research and pedagogic content that is increasingly abundant in digital open access format for text and data mining (TDM) purposes. This open access content is highly relevant to learning features of specialist varieties of English from across the academy but is otherwise off-limits for development into proprietary learning materials by the commercial education publishing industry (Fitzgerald et al., 2015, 2017; Wu et al., 2018). Indeed, the open corpus development work presented in this chapter would not have been possible had it not been for the campaigners for copyright reform, the Internet activists, the open policymakers, the open-source software developers and the advocates for open access, open data and open education that have made these resources available for re-use and remix.

This paradigm leads down several paths, including research into understanding how users actually perceive, appropriate and use the approach based on the open tools and resources provided. This inquiry informs their design and development in an R&D process that is presented here through the methodological lens of design-based research (Fitzgerald, 2019). This approach will be fundamentally different than if we

95 assume the user is actually a DDL or linguistics expert or that such an expert will be
96 the learner's interface to the system by preparing output for the learner to experience
97 and learn from (Johns, 1991b). This approach will necessarily also be different than if
98 we assume the user is always a formally registered student at a university with access
99 to EAP support that may or may not offer DDL or linguistics expertise for learning
100 the language features of specific discourse communities from across the academy.
101 The assumption behind this new paradigm that the right tools and resources can allow
102 the end-learner to drive the processes autonomously is fundamentally revolutionary.
103 This premise goes to the original contribution to the knowledge of this research but
104 also challenges and directs researchers and practitioners in the field to consider and
105 take up this new direction with open data-driven language learning systems design
106 for applications that can be scaled in higher education to meet the increasing numbers
107 of learners who are coming online in increasingly uncertain times (Fitzgerald, 2019;
108 König et al., 2022).

109 The focus on domain-specific language learning support via data-driven
110 approaches is, of course, also decidedly different from the current English for
111 Academic Purposes (EAP) paradigm, which in mainstream practice has been steadily
112 evolving away from its roots in English for Specific Purposes (ESP), domain
113 specificity and DDL processes towards the generic skills and knowledge programs
114 currently in vogue that are arguably being steered by generic EAP coursebook publi-
115 cations from the commercial education publishing industry (Gillett, 2018). Thus, this
116 is also a new paradigm based on DDL approaches, driving domain-specific language
117 learning support for EAP across formal, non-formal and informal learning modalities
118 in higher education. It will transform, potentially, the focus of DDL systems design
119 developments in language support and learning in general towards the non-specialist
120 end-learner but also hopefully help re-establish the centrality of language specificity
121 to the field of EAP (Anthony, 2018).

122 This new paradigm is necessarily rooted in greater multi- or trans-disciplinarity
123 (Colpaert, 2004, 2018). Given the goal of facilitating, in particular, the increasing
124 number of learners who are coming online in these uncertain times, and users of
125 large-scale MOOC platforms who are trying to function in domain-specific subject
126 areas that are invariably offered in the English language, the approach requires collab-
127 oration and cooperation among platform providers, subject academics and instruc-
128 tors, educational technologists, software developers, educational researchers, EAP
129 practitioners, linguists with expertise in corpus-based and DDL approaches and poli-
130 cymakers in knowledge organisations (libraries, universities, archives). It has to be
131 remarked, also, that the value and significance of this multi-disciplinary work is
132 amplified by our current situation in higher education with the pandemic, which
133 has seen a massive, urgent push to move learning online with an accompanying
134 impetus to identify, adapt and leverage learning content worldwide and to exploit
135 open educational resources, in particular.

6.2 Research Context

The open access movement in research and higher education has bolstered unprecedented access to artefacts of the academy in the form of published research articles, in addition to online platforms and services for accessing unpublished theses and pedagogic materials. One example is open access to transcribed video lectures and course reading content from the world's leading universities and institutions with an expanding provision in MOOCs. A further example is an open access to a growing corpus of over half a million PhD theses from universities across the UK with the British Library's Electronic Theses Online Service (EThOS). Both of these examples will feature for discussion in this chapter with respect to the nuanced meanings of openness and the tensions around human and machine re-use of content; the latter of which involves computational processes whereby texts and data are crawled and mined by software to build on and create new knowledge and derivative resources. Specifically, the research presented in this chapter is concerned with stakeholder reflections on a new paradigm for the co-design and co-development of data-driven language learning systems derived from open access content.

A definition for open access appeared for the first time in the declaration of the Budapest Open Access Initiative (BOAI):

By "open access" to [peer-reviewed research literature], we mean its free availability on the public internet, permitting any users to read, download, copy, distribute, print, search, or link to the full texts of these articles, crawl them for indexing, pass them as data to software, or use them for any other lawful purpose, without financial, legal, or technical barriers other than those inseparable from gaining access to the internet itself. The only constraint on reproduction and distribution, and the only role for copyright in this domain, should be to give authors control over the integrity of their work and the right to be properly acknowledged and cited (BOAI, 2002).

Colpaert (2016) from the field of Computer Assisted Language Learning (CALL) divides uses for data in language education into two main categories depending on divergent goals for re-use: data as content and data as information. The former data category includes authentic content found on the Web, including open access content that makes up the primary focus of this chapter. In contrast, the latter category includes information about data, otherwise known as metadata, which we also make use of in our research and refer to in this chapter.

Corpus linguistics researchers have demonstrated the importance of viewing language as data (Anthony, 2014; Hunston, 2002; McEnery et al., 2006; Sinclair, 2004). By way of extension, DDL can be viewed as a means for language teachers and learners to obtain, organise and study authentic language data derived from corpora in language education (Boulton & Cobb, 2017; Boulton & Pérez-Paredes, 2014; Boulton & Thomas, 2012; Chang, 2014; Cobb & Boulton, 2015; Vyatkina, 2016). There remains a persistent lack of exposure to and use of corpus-based systems and Natural Language Processing (NLP) tools by language practitioners in mainstream language education, however:

178 Many of the 15 million English teachers in the world today, according to the British Council
179 Annual Report (2010), have never heard of corpora, while many who are familiar with their
180 use by lexicographers and grammarians are not aware that they can use them themselves, as
181 could their students. (Thomas, 2017, p. 17)

182 For this R&D project, we identified a range of open, authentic domain-specific
183 text and data sources that are of perceived value to the EAP community yet are
184 off-limits for commercial re-use and development by the English language content
185 publishing industry. In this chapter, we will share reflections on our work with repre-
186 sentatives from knowledge organisations that manage and curate digital open access
187 content, such as the British Library, who are working at the cutting edge of reforms
188 in UK copyright law to create open access policy within their Research and Reuse
189 Committee. In line with the Fair Use Doctrine, which is a limitation to US copyright
190 law, an important exception and limitation to UK copyright law for TDM was intro-
191 duced in 2014, whereby permissions were established for the non-commercial re-use
192 of digital research content following an independent government report (Hargreaves,
193 2011).

194 One of the aims of this research has been to bring corpus linguistics researchers and
195 EAP practitioners to the interface of data-driven language learning systems design
196 for higher education through open initiatives in software development, research,
197 education and publishing that support the co-design, co-creation and distribution
198 of such systems. A further aim of this research has been to explore the potential of
199 working with open, authentic academic texts that afford language specificity (Hyland,
200 2002; Strevens, 1988) in the development of teaching and learning resources for
201 EAP that reflect the specific language and discourse features from target academic
202 communities.

203 We will discuss the perceived value that EAP researchers, teachers and managers
204 place on the efficacy of utilising authentic open access academic texts and corpora
205 in data-driven approaches for blended learning. These perceived educational values
206 will be weighed against the perceived risks held by knowledge organisations and
207 the individuals working therein, such as curators, subject academics and educational
208 technologists, regarding the remix and re-use of digital open access content and
209 collections for non-commercial research and education purposes. For the scope of
210 this chapter, we will explore the following research questions:

- 211 (1) To what extent can open access content foster open educational practices
212 among academic English language stakeholders for designing, developing and
213 evaluating data-driven language learning resources?
- 214 (2) What impact do the underlying business models and cultural practices of insti-
215 tutions and organisations have on open educational practices for remixing open
216 access content in the design, development, implementation and dissemination
217 of resources for EAP in higher education?

218 6.2.1 Research Materials

219 With this research, we have placed particular emphasis on co-designing and co-
 220 creating language learning systems for pedagogic purposes rather than for corpus
 221 linguistics research purposes. Drawing on the concept of knowledge mobilisation
 222 (Levin, 2011), our goal is to engage relevant stakeholders in moving available knowl-
 223 edge from research in corpus linguistics, open education and computer science (NLP
 224 and TDM) towards knowledge users, namely EAP practitioners and learners. The
 225 goal is for knowledge users to not only benefit from the research but to collabo-
 226 rate directly in an iterative design-based research process. Intermediaries working
 227 in knowledge organisations have acted as brokers and open education champions in
 228 this research by creating access to knowledge artefacts that are valued for re-use in
 229 EAP via initiatives in open access policy and reforms in copyright law.

230 Although the findings from this research are tied to issues with designing and
 231 developing open access content into data-driven learning systems, wider issues vis-
 232 à-vis the re-use and remix of open access content in language materials development
 233 practices will also be discussed as they apply to both classroom teaching and online
 234 learning. The Appendix at the end of this chapter provides an overview of our work
 235 to date. It identifies the knowledge organisations, researchers and knowledge users
 236 who have collaborated on the design and development of open data-driven systems
 237 for learning aspects of academic English in formal and non-formal higher education
 238 contexts with the FLAX¹ and F-Lingo² projects.

239 6.3 Research Methods and Results

240 Methods for collecting data from different participant groups in different loca-
 241 tions over a period of years included: focus discussions, face-2-face and Skype
 242 interviews and email exchanges stemming from project meetings on observations
 243 and evaluations shared in this situated research. Three knowledge organisations
 244 have participated in the research (The British Library, The Oxford Text Archive
 245 and the Connected Repositories research group at the UK Open University). Eight
 246 researchers working in the area of corpus and computational linguistics and open
 247 education have participated in the research from higher education institutions in
 248 Aotearoa, New Zealand, Spain, Canada and the United Kingdom. Seven knowledge
 249 users working in EAP teaching and management from two UK universities have
 250 also participated in the research. Automated content analysis (ACA) was carried out
 251 on the complete corpus employing the Leximancer software version 4.5, and then
 252 on sub-corpora corresponding to data from the three stakeholder groups engaged in
 253 this research—knowledge organisations, researchers and knowledge users. Results
 254 from the ACA in this study were checked and then triangulated with participants

¹ <http://flax.nzdl.org/greenstone3/flax>.

² <https://chrome.google.com/webstore/search/flingo>.

255 in this qualitative research to create opportunities for participants to comment on
 256 transcripts and emerging findings. Thematic and conceptual findings in the datasets
 257 were then confirmed with participants as they pertain to reflections on the itera-
 258 tive design processes for designing open data-driven systems for academic English.
 259 The complete corpus and ACA visualisation maps of key themes and concepts from
 260 this R&D project are available for viewing on the Open Science Foundation³ data
 261 platform.

262 **6.3.1 Design-Based Research**

263 Action research is a widely employed methodology in English language education
 264 research and teacher training programmes (Burns, 2009) and shares many of the same
 265 principles as design-based research (DBR). Pragmatism is central to both approaches,
 266 often employing mixed methods of inquiry to arrive at tangible solutions to educa-
 267 tional problems. Within action research cycles, individual teaching practitioners carry
 268 out classroom teaching interventions to observe, record and reflect on the impact of
 269 these interventions over time to inform and improve their classroom and online
 270 teaching practice (Reason & Bradbury, 2007). In design-based research, another
 271 layer exists that requires educational practitioners to collaborate with research and
 272 design teams (Anderson & Shattuck, 2012).

273 Although DBR has sustained great interest from researchers and practitioners
 274 within the instructional design and educational technology milieu, it is nevertheless
 275 a long-term and very resource-intensive exploratory research method with goals
 276 and outcomes that are difficult to define. The literature on DBR attests to “a series
 277 of approaches, with the intent of producing new theories, artefacts, and practices”
 278 (Barab & Squire, 2004, p. 2). More specifically, these approaches have been defined
 279 as multiple research cycles that include numerous iterations of analysis, design,
 280 development, evaluation and revision (Burkhardt, 2006; Walker, 2006; Amiel &
 281 Reeves, 2008; Hakkarainen, 2009; McKenney & Reeves, 2012). Data are collected
 282 over a minimum of several weeks but, in most cases, are collected over several months
 283 or years (Herrington et al., 2007) as has been the case with our research, which has
 284 been ongoing for over a decade now (Fitzgerald, 2019).

285 **6.3.2 Automated Content Analysis**

286 Automated Content Analysis (ACA) is situated within the framework of compu-
 287 tational social sciences. It refers to a range of algorithms that employ probabilistic
 288 models, namely topic models and concept mapping models (Blei, 2012a, b), that iter-
 289 atively infer the themes and concepts present within a corpus. ACA can be traced back

³ <https://osf.io/gbkzp/>.

290 to the theoretical underpinnings of Latent Semantic Indexing (LSI; Papadimitriou
291 et al., 1998), leading to the three-level Bayesian model of Latent Dirichlet Allocation
292 (LDA; Blei et al., 2003). The current state of the art with ACA models involves the
293 identification and analysis of higher levels of complexity found within thematic struc-
294 tures (Blei, 2012b). Current ACA systems include features for analysing “syntax,
295 concept hierarchies, document networks and temporal trends in themes, furthering
296 our ability to visualize and explore the literature” (Nunez-Mir et al., 2016). ACA is
297 primarily used to automatically analyse text in digital format but also, increasingly,
298 media content, e.g., images (Boumans & Trilling, 2016).

299 In this section, we look through the analytical lens offered by ACA at the different
300 themes and concepts from each of the three participant groups in this study: knowl-
301 edge organisations, researchers and knowledge users. Due to the limited scope of
302 this publication, we will only be looking at the results of the top four themes in
303 each sub-dataset for the three participant groups in this research. Where we present a
304 summary and discussion of results from all three sub-datasets, themes and concepts
305 from the data will be italicised.

306 Our reasons for employing the Leximancer ACA software to analyse the quali-
307 tative datasets were two-fold: to increase validity and to determine the lexical co-
308 occurrence of information extracted from natural language into semantic or concep-
309 tual patterns using automated methods. Leximancer has been designed to miti-
310 gate subjectivity and researcher bias in the traditional content analysis processes
311 of manual text analysis, coding and intercoder reliability testing (Weber, 1990).
312 Through powerful automated methods, Leximancer is devised to make the human
313 analyst aware of “the global context and significance of concepts and to help avoid
314 fixation on particular anecdotal evidence” (Smith & Humphreys, 2006, p. 262). Lexi-
315 mancer performs two types of analysis on a ranked list of lexical terms found in
316 a unified body of text or corpus: conceptual analysis and relational analysis. The
317 conceptual analysis measures the presence and frequency of concepts in a document
318 set by extracting words, phrases, or collections of words that represent a concept. The
319 relational analysis is concerned with measuring the co-occurrence of concepts within
320 a document set, extracting these co-occurring concepts to show their relationship.

321 The design principles that underpin the Leximancer software are founded on
322 observations from the fields of corpus linguistics, computational linguistics and
323 psycholinguistics, resulting in the development of the semantic and relational Lexi-
324 mancer algorithms that are employed in both stages of the software’s co-occurrence
325 information extraction technique (see Smith, 2000a, b, 2003). Leximancer was
326 employed to mine the total qualitative dataset and sub-datasets for each partici-
327 pant group, resulting in a thesaurus of words identified within each corpus analysed
328 along with their related meanings and surrounding words or collocates. The complete
329 corpus and ACA visualisation maps of key themes and concepts from this R&D
330 project are available for viewing on the Open Science Foundation⁴ data platform.

331 As shown in Fig. 6.1, closely related words from the complete qualitative dataset in
332 this study are identified by the ACA software as concepts and are represented as dots

⁴ <https://osf.io/gbkzp/>.

333 within thematic circles of interrelated concepts on a concept map. The key below the
 334 map indicates how many times the central themes occurred in the corpus. Important
 335 themes are mapped with warm colours. For example, *research* and FLAX appear
 336 in red and brown on the concept map (Angus et al., 2013). These two dominant
 337 themes are represented as tightly packed circles containing concept dots in close
 338 proximity to one another. The spatial alignment of these dots indicates how closely
 339 related concepts are within each key theme (Smith & Humphreys, 2006). For instance,
 340 *research, corpus, able, EAP, teaching* and *learning* are closely related concepts within
 341 the dominant *research* theme. Thematic circles are sometimes shown as overlapping
 342 with one another when concepts occur close to or across neighbouring themes, such
 343 as the concepts for *corpus* and *learning* within the *open* and *research* themes, which
 344 are central to this ongoing design-based research with the FLAX project and will
 345 provide a basis for the discussion section of this chapter.

346 6.3.2.1 Knowledge Organisations

347 The Leximancer analysis of data from the knowledge organisations group reveals
 348 *text* as the major theme. The concepts within this key theme of *text* emphasise *exper-*
 349 *imentation* with *corpora* and *stuff*, with one frequent example in the dataset being
 350 the *EThOS* (Electronic Thesis Online Service) PhD thesis content at the British
 351 Library, in addition to the *terms* around *re-use*, and what you are *able* to do when
 352 *using* texts with *text* and *data mining*. The second most prominent theme is *work*
 353 with concepts reflecting the importance of *doing work* in the *open* as central to this
 354 design-based research with knowledge organisations. In close orbit to the *text* theme
 355 are the overlapping themes of *trying* and *example*, representing the third and fourth
 356 most frequent themes in the dataset, coming in closely behind the *work* theme. Of
 357 note in the *trying* theme are the connected concepts of *people trying* to do *things*.
 358 *Re-use* is the concept shared between the overlapping *text* and *example* themes. Also
 359 apparent in the *theme are the key interlinked concepts of example, collections and*
 360 *metadata for what can probably be looked at with respect* to research and develop-
 361 ment that focus on the *re-use of text* and their *metadata* from digital *collections*. In
 362 the discussion section, we will explore these themes and concepts further with refer-
 363 ence to the terms and conditions around open access content re-use in this research
 364 with knowledge organisations.

365 6.3.2.2 Researchers

366 We now turn to interview data with education researchers who have worked with the
 367 FLAX project. The first researcher interviewed was Researcher 4, a legal English
 368 corpus researcher at the University of Murcia in Spain who developed the British
 369 Law Reports Corpus (BLaRC) with judicial hearings from around the world that
 370 subscribe to the English common law system. The corpus was made available with
 371 an open access government licence from the British and Irish Legal Institute (BAILI).

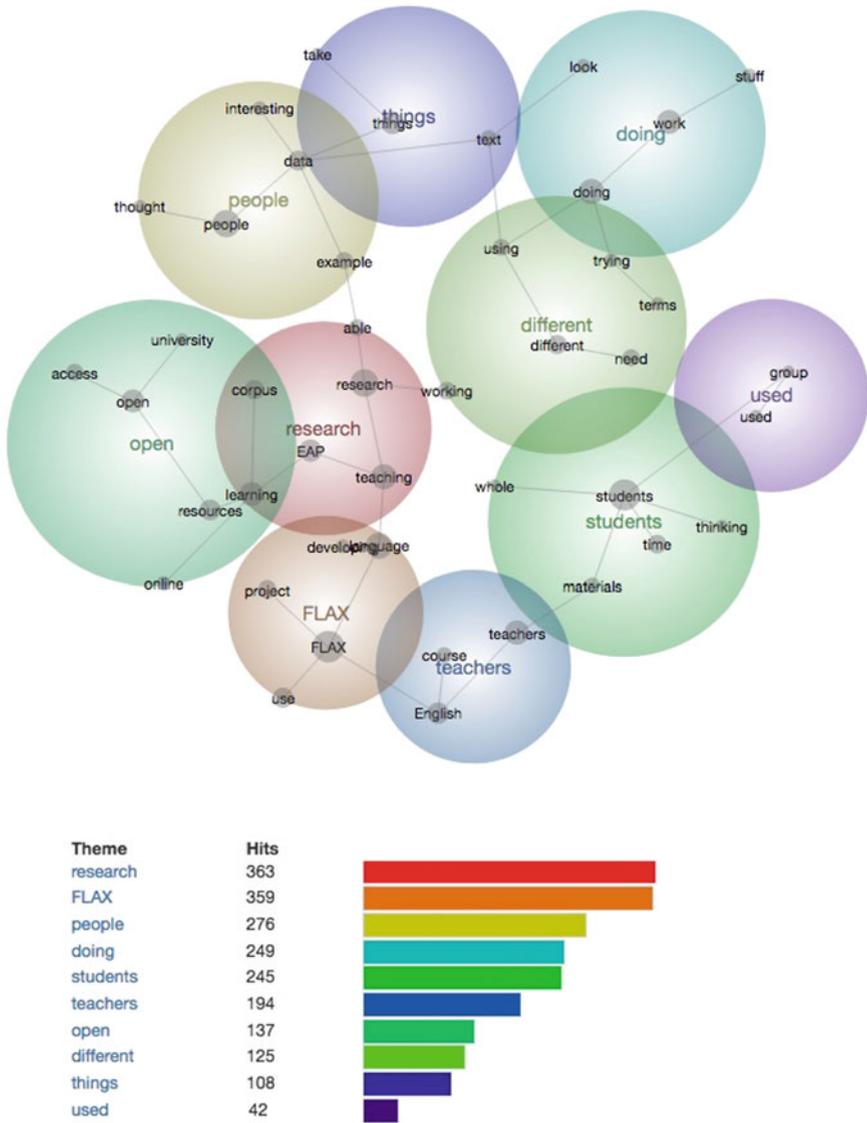


Fig. 6.1 Concept map and key derived from automated content analysis of the complete qualitative dataset

UNCC

372 Researcher 6 also conducted doctoral research into lexical bundles with the FLAX
373 project, focusing on the Chinese and New Zealand EAP contexts.

374 When we look at the Leximancer conceptual analysis for the researchers' group, of
375 note are four prominent and overlapping themes: *FLAX*, *students*, *teachers* and *time*.
376 Moreover, the concepts of *access*, *different*, *research*, *online*, *language* and *learning*
377 appear in the overlapping foci areas of these top four central themes. In this section,
378 we will summarise the findings from these concepts, which will form the basis for
379 the discussion section of this researcher participant group later in the chapter. The
380 *access* concept in particular, which appears in the overlap between the *FLAX* and
381 *students* themes in the ACA, is expressed in the data as issues related to conducting
382 *research* that provides students with *access* to and *use* of *different corpora*, *data*
383 and *systems* in *FLAX* that can support their *online language learning* with formal
384 *language courses* and non-formal *MOOCs*. Of interest, the *access* concept is also
385 expressed in the data in relation to the issue of *gaining access* to *students* through
386 *working with language teachers* to conduct *research* into the *use* of the *FLAX system*.
387 This last point on *access* is further extended into the sixth most frequent theme in the
388 dataset, *study*, with concepts expressing the *need* for *user studies* on the uptake of
389 *FLAX*. In addition, the issue of *access* is further expressed by how *teachers* may be
390 *interested* in *working* with the *FLAX project* but are limited in terms of the fourth most
391 frequent theme, *time*, due to the heavy emphasis placed on *teaching* and *learning*
392 and not on conducting *research* at their institutions.

393 6.3.2.3 Knowledge Users

394 Of the seven EAP practitioners who participated in the research, only one (Knowledge
395 User 1) from Queen Mary University of London (hereafter referred to as QMUL),
396 had extensive experience with using corpus tools in his classroom teaching, namely
397 the Sketch Engine⁵ suite of tools for querying and sketching corpora. The three other
398 participants at QMUL (Knowledge User 2, Knowledge User 3 and Knowledge User
399 4) all had a background in CALL for developing online EAP resources, most notably
400 Academic English Online.⁶ The three EAP teachers at Durham University who are
401 former EAP teaching colleagues of Researcher 1 (Knowledge User 5, Knowledge
402 User 6 and Knowledge User 7), were early adopters and advocates for using open-
403 source software and OERs in their classroom teaching as a means of ensuring that
404 their students had access to high-quality free and open online teaching and learning
405 resources during and after their courses had finished. The EAP practitioners in this
406 study expressed that the motivation to adopt open educational practices as they apply
407 to academic practice in higher education was a motivating factor for participating in
408 the research with the FLAX project.

409 The dominant themes arising from the Leximancer analysis of interviews and
410 focus discussions from project meetings with knowledge users—EAP teachers and

⁵ <https://www.sketchengine.eu/>.

⁶ <http://aeo.sllf.qmul.ac.uk/>.

411 course managers—are *EAP* followed closely by *students*, *things* and *people*. In
 412 summary, results from the ACA of this sub-dataset point to issues concerned with
 413 the concepts of *EAP* and the *teaching* of *academic English language* from the largest
 414 theme, *EAP*. The second-largest theme in the data, *students*, reveals issues around
 415 *materials* for *teaching students* that *teachers* are developing themselves or those
 416 *materials* that have been developed by commercial publishers and reflections on what
 417 does and does not *work* in practice. The third most frequent theme in the dataset,
 418 *things*, is representative of concepts related to what *needs* to be done with *research*
 419 using *things* and *materials*. In the fourth most frequent theme, *people*, an interesting
 420 interplay of concepts are revealed in reference to *people* as being those *EAP* teachers
 421 working in universities who do or do not create *access* to *open resources* for *educa-*
 422 *tion*, and also in reference to *people* outside of the university who can and cannot
 423 *access open resources* for education. The themes and concepts outlined here in this
 424 section will be explored in more depth in the corresponding discussion section of
 425 this paper on knowledge users.

426 The work at Durham University took the form of an OER cascade training project
 427 with the participating *EAP* practitioners and their students that introduced them to
 428 four online data-driven text analysis language learning systems: Lextutor,⁷ AntConc,⁸
 429 Word and Phrase⁹ and FLAX. This OER cascade training work led to collaborative
 430 evaluations and further development iterations of the Learning Collocations collec-
 431 tion in FLAX with the addition of the open access British Academic Written English
 432 (BAWE) corpus managed by the Oxford Text Archive for a specific focus on academic
 433 English collocations. It was written up as a case study for the UK Higher Education
 434 Academy (Fitzgerald, 2013). This work at Durham also resulted in the development
 435 of the full-text BAWE collections in FLAX that focused on novel ways to search and
 436 browse augmented academic texts that represented different genre types from across
 437 the disciplines of the arts and humanities, the social sciences, the physical sciences
 438 and the life sciences (Wu & Witten, 2016).

439 The work at QMUL focused on design collaborations with open access PhD thesis
 440 abstract content managed by the British Library for the development of domain-
 441 specific micro-corpora and interactive games with Android mobile apps for uptake
 442 on QMUL's pre-sessional *EAP* programmes (Fitzgerald et al., 2014). The work with
 443 QMUL led to a further design iteration with the development of the much larger PhD
 444 Abstract collections in FLAX of 9.8 million words (Wu et al., 2018).

445 6.4 Discussion

446 This section discusses prominent themes and interrelated concepts from the ACA of
 447 the qualitative datasets. We drill further down into the data that captured reflections

⁷ <https://www.lexutor.ca/>.

⁸ <http://www.laurenceanthony.net/software.html>.

⁹ <https://www.wordandphrase.info/>.

448 from participants in the research to present relevant themes and concepts identified
449 in transcriptions. Where we refer to actual data for discussion, themes and concepts
450 will be italicised.

451 **6.4.1 *The Four Pillars of Re-Use in Knowledge Organisations***

452 Our research with knowledge organisations in developing open corpora for EAP
453 shows that it often comes down to those individuals working on the inside who
454 are reasonably au fait with copyright law as it pertains to open access and open
455 educational practices and who are willing to champion the re-use of resources and
456 encourage the development of open policies within their organisations. We have seen
457 this type of open access policy championship with the EThOS service team manager
458 and the British Library Labs project manager. The progress with policy development
459 for open access and re-use that enable TDM approaches with digital collections at
460 public knowledge organisations such as the British Library is contrasted with the
461 absence of open education policy in higher education institutions, where there has
462 been less progress made with the re-use of educational content. Open access, in most
463 cases, to read-only research publications and, in lesser cases, to pedagogic content
464 has become the default re-use position of most universities and mainstream MOOC
465 providers.

466 Once again, those individuals who are already open education practitioners who
467 have openly licensed their educational resources with Creative Commons licenses
468 have enabled the FLAX team to develop derivative language learning collections.
469 Open licensing supports their wider practices in open digital scholarship (Weller,
470 2011)—via blogs, public lectures, MOOCs, networked courses etcetera—to widely
471 promote the subjects they are passionate about. Notably, Professor Fisher of the
472 CopyrightX micro-networked course has deliberately applied his expertise in under-
473 standing the ins and outs of copyright law by licensing his teaching and learning
474 content as CC-BY with Creative Commons “to maximise the number and variety of
475 educational projects and derivative works that can be built (directly or indirectly) on
476 our foundation—and thus the set of students who might benefit from our efforts”.
477 (Fisher, 2014, p. 17).

478 The participating knowledge organisations in this research differ regarding poli-
479 cies and practices around re-use. British Library Labs (BL Labs) is an Andrew
480 Mellon Foundation-funded initiative that supports the remix and re-use of the British
481 Library’s digital collections and data for research and educational purposes. In an
482 interview with the project manager of BL Labs (Knowledge Organisation Repre-
483 sentative 1 in the transcript corpus), we discussed the FLAX project research with
484 the EThOS PhD theses dataset for the development of the PhD Abstract collections,
485 wherein he identified four pillars for the re-use of this dataset that can be broadly
486 applied to the re-use of other digital collections at the British Library:

- 487 (1) “Do we have an expert with curatorial knowledge of a particular *collection* who
 488 is on board with *re-use*? Some curators are not concerned about that at all. All
 489 they care about is the preservation and not about who *uses* it.
 490 (2) Do we know where it, the collection, is? A description of something is one thing
 491 but who actually has the digital files? Can they be accessed?
 492 (3) Is there any *metadata*? That obviously helps enormously because it means that
 493 you can then release the *metadata*, normally. But even *metadata* has licenses as
 494 well... so, who owns that *metadata*?
 495 (4) Is the *collection* close to being copyright-cleared? And what I mean by that,
 496 I actually mean, is it, could it potentially, easily, be available under an *open*
 497 licence?”

498 (Interview with BL Labs Manager, British Library, London, UK)

499 With the harvested PhD theses in EThOS at the British Library, the provenance
 500 is very mixed, whereby there is no one set of terms and conditions for re-use of the
 501 open access content found therein. This phenomenon is largely a reflection of the
 502 different universities where the research was carried out and is dependent on whether
 503 or not there were industry investments in the research, for example, which would
 504 result in copyright stakes. Due to this mixed provenance, the British Library has
 505 undertaken measures to balance any possible research instances of re-use with any
 506 identifiable potential risks such as mass copying, misrepresenting and misquoting of
 507 the EThOS dataset. As with the Oxford Text Archive, a cautious approach has been
 508 adopted at the British Library with respect to TDM, whereby collections are only
 509 available for non-commercial re-use purposes on a request-only basis. The BL Labs
 510 manager does, however, acknowledge the iterative nature of research and encourages
 511 the practice of “dogfooding” at the British Library, whereby collections management
 512 teams, such as the EThOS team, engage in internal research on the re-use and remix
 513 of collections to anticipate affordances and hindrances with conducting research:

514 Knowledge Organisation Rep 1: First, to *work* with a collection it’s important to ensure that
 515 there’s a human being who can tell you the story of that collection because you don’t know
 516 what may be lurking in there and it may not be about legal issues. It could be political. It
 517 could be financial. But that information isn’t always documented.

518 Researcher 1: Sorry to interrupt you there, but were there any issues around *EThOS*?

519 Knowledge Organisation Rep 1: Well, I think there are still issues really because the problem
 520 of *doing* this work is because the intellectual property is going to be dependent on the
 521 institution and their relationship with their students. It seems that that is not straightforward
 522 with all the different institutions. So, if you do a PhD at an institution, you’re under the IPR
 523 for that *work*, and I think that different universities have different views and policies.

524 Researcher 1: So, it’s not always automatically the student’s *work*? I thought it was.

525 Knowledge Organisation Rep 1: All I know is that some *work*, some PhD *work*, is embargoed
 526 because it has commercial sensitivities in there. So, for *example*, somebody might...

527 Researcher 1: Because they’ve been funded by...?

528 Knowledge Organisation Rep 1: Yeah, because they've been funded by Panasonic, for
529 *example*.

530 Researcher 1: Yeah, I get that.

531 Knowledge Organisation Rep 1: There could be, depending on the PhD and the funding
532 stream, so it could not only be the university, it could be the funder, the funder might have
533 certain requirements. It could be commercial. It could be a funding council. What you're
534 getting is a harvested bunch of stuff in *ETHOS* where the provenance is very mixed, and I
535 think the team have decided to take a very cautious approach in terms of being *able* to do
536 things like *text* and *text* and *data mining*, so, you know, it's on a request only basis. Because,
537 especially, you know, about the possibility that there could be commercial *re-use*.

538 Researcher 1: Yes, I think that's getting back to your original point about the library wanting
539 to know what your research questions were before *doing* the *work*.

540 Knowledge Organisation Rep 1: Exactly.

541 Researcher 1: And that's when somebody puts in a request, for *example*. We want to *re-use*
542 these *texts* for these purposes, and this is what the end result will *look* like kind of.

543 Knowledge Organisation Rep 1: Yeah, but the problem with that is, in our experience, is that
544 research doesn't *work* like that. With research you don't know what you're going to get. You
545 might know your research questions, but the whole point and nature of research is that it's
546 iterative. You know, you *experiment*.

547 Researcher 1: I'm glad to hear you say that because, you know, that was our experience
548 with the Oxford Text Archive when we requested the BAWE *corpus*. Because we didn't
549 know in advance that we'd be Wikifying whole *texts* but then we had the technology to do
550 it. In particular, I mean, all the prior *work* we had done with Wikipedia *mining* at the Digital
551 Library Lab at Waikato. And, we thought, well, Wikification may well be useful for language
552 learning so let's add this functionality for learners. So, the BAWE *collections* became our first
553 Wikified *collections*, and you can see this feature in our subsequent *collections*, including
554 the PhD Abstract *collections* with *ETHOS metadata*. But this *work* with Wikification wasn't
555 in our initial request to the OTA, which was instead very general in terms of what we were
556 proposing to do.

557 Knowledge Organisation Rep 1: Yeah, I think in general, I understand why there needs to
558 be this clarity but unfortunately, it's a complete misunderstanding of the whole scholarly
559 process. The scholarly process is actually incredibly creative, and you know, you don't know
560 by the very nature of research, that you don't know what you're going to find. And, you know,
561 it's surprising what comes along the way. Ideas will come along the way, and that's just the
562 nature of research. So, we have found that really challenging. And, what we've decided to
563 do, I think, is to be *working* on research questions where they can be sort of dealt with on a
564 case-by-case basis, and also to agree on what the outcomes are going to be. So that, like, if
565 *people* want to publish *work*, what actually can be published, and what can't be published
566 because of the sensitivities at the moment. We're also having quite a lot of requests to do
567 *text* and *data mining work* with our non-print legal deposit *stuff*.

568 (Interview excerpt with BL Labs Manager, British Library, London, UK)

6.4.2 Issues of Access in DDL Research

The automated content analysis of the entire qualitative dataset reveals a direct link between the knowledge organisations and researchers' sub-groups with overlapping themes of *access*. Put simply, access to digital collections that can be re-used by researchers, in this case, corpus linguistics and open education researchers, is due in no small part to the open access and open education policies adopted by knowledge organisations and the gatekeepers working within those organisations who implement these policies to promote open access and re-use.

We turn first to a discussion with Researcher 4 in this study on the perceived affordances of re-using and remixing open access publications for open data-driven learning in DDL research with reference to the BLARC of 8.85 million words (Marín et al., 2014), which is derived from open access judicial hearings licensed with a government license and available from the BAILII online service. Marín developed the BLARC due to the lack of relevant, authentic materials for teaching the specific area of legal English in EAP. We invited her to include her corpus on the FLAX website so that it would be openly accessible for data-driven language learning in addition to corpus linguistics research. Researcher 4 was interviewed about the making of the BLARC, which highlights the affordance of the *access* concept as a prominent concept in the interview data with applied corpus linguistics researchers, and how this had enabled the development of legal English resources from open access content in comparison with proprietary legal content services that require licence subscriptions:

Researcher 1: You know, my next question: Could you even have built the BLARC without those open government licenses on all of those documents, those judicial hearings in the BAILII (British and Irish Legal Information Institute)?

Researcher 4: No, that's the thing, that's the thing. The amazing discovery was the BAILII [...] I was *thinking* about buying a licence for LexisNexis, I think it's called. There are a couple of them, which cost a fortune, a fortune. I'm not sure but I think law firms, they pay, I don't know, four or five thousand pounds a year for having that kind of *thing*, which is amazing [...]

Researcher 4: Actually, the University of Murcia doesn't have *access* to that database because one of my colleagues was in Madrid, she was a visiting researcher there, and she downloaded like a hundred thousand texts from LexisNexis because she didn't know that the BAILII existed. So, when she came here, and we were talking, and I said, *look* there's this site [the BAILII] and they have added a lot of overseas legal documents, including United States documents. They have the whole planet in there. It's amazing how much stuff you can find. So, to me it was a huge, huge discovery. That was the best thing that could have happened to me. That's why I started my *research* on legal *corpora*. I mean that was one of the reasons.

Researcher 1: *Access* is so key, isn't it? And I'm sure that's a big part of why the BAILII exists as well because they knew people couldn't *access* LexisNexis.

(Interview excerpt with Researcher 4, via Skype from Murcia, Spain)

The experience of Researcher 6 in this study regarding attempts to carry out DDL research with language teachers and learners in China highlights another aspect of

612 the *access* concept as it intercepts with the dominant themes for *FLAX*, *students* and
 613 *teachers* within the qualitative dataset. Her greatest challenges were with securing
 614 *access* to research sites with *students* and *teachers* in China to test out the efficacy
 615 of the Learning Collocations collection in the *FLAX* system. She and Researcher 1,
 616 both of whom come from the field of education, discussed the role of *use* or user
 617 *studies*—prevalent concepts within the data—with tools and projects like *FLAX* that
 618 stem from computer science as they are applied to the *students* theme for educational
 619 researchers:

620 Researcher 1: They talk a lot about *user studies* in computer science, don't they?

621 Researcher 6: Yeah, but those *user studies* are only to prove that the tool works.

622 Researcher 1: Right, the focus is not to prove that *learning* has occurred with use of the tool.

623 Researcher 6: No, the purpose of such *user studies* in computer science is not to promote
 624 the application of the tool. So, for them the end of their *project* is that the tool has been
 625 developed successfully but for *English teachers* with *English language learning* tools, that
 626 is the beginning. But between the end of computer scientists completing the development of
 627 a *learning* tool and the beginning of *English language teachers* adopting a *learning* tool in
 628 their *teaching* there is a gap.

629 (Interview excerpt with Researcher 6, University of Waikato NZ)

630 The importance of user studies in this design-based research leads into our final
 631 section of analysis on the data collected with knowledge users, EAP teachers and
 632 managers at two UK universities, Durham and Queen Mary.

633 **6.4.3 Barriers to Remixing Texts in Data-Driven EAP** 634 **Materials Development**

635 Collaborative work with Durham and Queen Mary revealed that data-driven
 636 approaches are not embedded within materials development and classroom teaching
 637 practices at these two UK universities. However, online corpus-based resources have
 638 a valued place as supplementary EAP materials. Most DDL tools and corpus-based
 639 systems were viewed by the majority of participants at Durham and QMUL as stand-
 640 alone web-based reference resources for students to explore outside of classroom
 641 teaching time.

642 Issues stemming from the design-based research carried out with Durham and
 643 QMUL include the limited amount of time EAP teachers have in the classroom with
 644 students to focus on discrete language items and the infeasibility of shepherding large
 645 groups of students in developing and mining personalised domain-specific corpora
 646 for focused help with dissertation and thesis writing, for example. This is despite
 647 some promising findings from research into DDL approaches with smaller, more
 648 tailored EAP classes for building Do-It-Yourself digital corpora with students to
 649 help with PhD thesis writing (Charles, 2012, 2015).

650 The focus-group discussions with managers at QMUL on the increased availability
 651 of open access content point to what EAP practitioners are now *able* to do with
 652 *academic things, resources* and *materials* for *use/using* with *students* as they emerge
 653 in this sub-dataset for the top four themes related to knowledge users: *EAP, students,*
 654 *things* and *people*. Knowledge User 2, manager of multimedia language support at
 655 QMUL, describes the approach of developing transferable skills in EAP materials
 656 development with revising and repurposing open access research publications as
 657 being one that is closer to traditional approaches with the re-use of authentic language
 658 content for classroom teaching purposes:

659 Knowledge User 2: You know, I think the thing about *open* educational *resources*, the question
 660 here, or part of the question here, which we discovered in this *project*, for example, is if you
 661 *take* a text, a raw text, which is not adapted for *teaching* like an article, it has *EAP* potential
 662 because it's an authentic *academic* article. Then the *ability* to use that and to put it into
 663 *materials*, or adapt it, modify it, or change it under the Creative Commons thing is the
 664 revelation. Because we've all been doing it for years anyway, from copying it from a book
 665 or something when we've not supposed to have been adapting it, changing it, or whatever.

666 (Knowledge User 2, focus-group discussion excerpt, Queen Mary University of London UK)

667 From the same focus discussion, the pre-sessional course director at QMUL,
 668 Knowledge User 3, talks about the barriers to people working in universities from
 669 openly sharing EAP materials across institutions and how they are tied to each univer-
 670 sity's business model with the aim of promoting their particular brand of EAP courses
 671 and materials as a unique selling feature. He also discusses the rise in influence of
 672 commercially produced EAP publications and the re-use of third-party materials
 673 from these publications as seeping into university EAP course materials develop-
 674 ment practices, which in turn creates a further barrier to sharing due to copyright
 675 infringement:

676 Knowledge User 3: There is a certain degree of *openness* but there is also this desire for
 677 everything to be branded, and a certain amount of clutching to your chest, especially about
 678 pre-sessional *materials*. [...] This is Queen Mary *material*, this is Southampton *material*,
 679 this is Durham *material*. But I think when you get back to the institutional level, those are
 680 where the real barriers lie because *people* are, and that comes down to the cut n paste culture
 681 that means a lot of third-party *materials* end up in our *materials* and are branded as being
 682 in-house but a lot of them are not really. You know, the ideas come from published *materials*
 683 and they're probably not properly acknowledged anyway because they're only being used
 684 internally. And part of that barrier to sharing more *openly* is raising an awareness of our
 685 existing practices and this means they don't want to share between institutions because
 686 they're worried that *people* will see just how much cut n paste is going into those *materials*.
 687 And I think the loser is the *student*, you know, because if *people* were really producing and
 688 sharing the best that they could amongst institutions to then create the best *EAP* pre-sessionals
 689 then the *students* would obviously benefit.

690 (Knowledge User 3, focus-group discussion excerpt, Queen Mary University of London UK)

691 From a meeting with QMUL EAP teacher, Knowledge User 1, the concepts of
 692 *open* and *access*, which congregate in the *people* theme, relate to frequent references
 693 in the data of how people outside the university can also benefit from *education* and
 694 *resources* that are openly accessible via the Internet:

695 Knowledge User 1: This *open*-source software and *open access* approach to data-driven
 696 *learning resources* does threaten current business models in *EAP* provision, doesn't it?
 697 This idea of yours to re-use the artefacts of the academy. This really bucks some *people* in
 698 academia.

699 Researcher 1: Tell me more about that because that's what I think is important to be doing
 700 in higher education, but I realise that this isn't everyone's priority.

701 Knowledge User 1: That's what I think is important as well. It's the ivory tower, isn't it? It's
 702 the secret garden behind the firewall of the ivory tower.

703 [...]

704 Knowledge User 1: Now, yes, I need *people* within this higher education environment [Queen
 705 Mary] to re-use these *academic* texts but I also *need people* to come into this *FLAX* environ-
 706 ment, *people who need* to interface with this environment for whatever *academic English*
 707 *need* they have, and that's what *FLAX* does for them in a manageable way. It makes it
 708 *accessible* not only to *people* who are using it in situ within the privileged brick-n-mortar of
 709 the academy but for *people* who, like I say, *need* to interface with that in some way outside of
 710 the academy, and, oh, that matters. The *resource* is not just locked inside our intranet-based
 711 VLE [Virtual Learning Environment] where I have developed *learning resources* with links
 712 out to *FLAX* on the web, which is really a Mickey Mouse version of *FLAX* in here.

713 (Meeting excerpt with Knowledge User 1, Cutty Sark pub in Greenwich, London, UK)

714 6.4.4 A Crisis in EAP Identity

715 An emerging tension in formal EAP is the issue of EAP practitioner identity in the
 716 neoliberal university (Ding & Bruce, 2017; Hadley, 2015; Hyland, 2002). Where
 717 are EAP service units placed in universities, and more importantly, how are they
 718 received and perceived by the wider academy? At its best, EAP is viewed as drawing
 719 on and contributing to a rich knowledge base from research in systemic functional
 720 linguistics, genre theory, corpus linguistics, academic literacies and critical EAP
 721 (Ding & Bruce, 2017). At its worst, EAP has been conceived as having "accepted
 722 the role as an economic and intellectual short-cut... [with] maximum throughput of
 723 students with minimum attainment levels in the language in the shortest possible
 724 time". (Turner, 2004, pp. 96–97).

725 There has been an upswing in commercially produced EAP publications with a
 726 notable shift in focus towards generic academic skills and processes. The increasing
 727 prominence of generic EAP publications can be seen to exacerbate the growing fissure
 728 in EAP practitioner identity with the emergence of two opposing camps: English for
 729 General Academic Purposes (EGAP) versus English for Specific Academic Purposes
 730 (ESAP). The received definitions and understandings from the literature indicate that
 731 EAP is a subset of English for Specific Purposes (ESP) (see ETIC, 1975; Widdowson,
 732 1983; Swales, 1985; Flowerdew & Peacock, 2001; Howatt, 2004; Belcher, 2010;
 733 Charles & Pecorari, 2016; Anthony, 2018). However, this understanding of EAP as
 734 being concerned with the teaching and learning of domain-specific language appears
 735 to have become conflated and confused as the popularity of generic skills-based

736 EAP textbooks, subscription-based supplementary online resources and programmes
737 continues to rise (Gillett, 2018).

738 The absence of data-driven approaches in the design of EAP classroom teaching
739 and online materials is a recurring theme in the sub-dataset from knowledge users.
740 In a focus-group discussion with former teaching colleagues at Durham (Knowledge
741 User 5 and Knowledge User 6), reflections turned towards collaborative work that
742 involved trialling corpora and data-driven approaches for EAP (Fitzgerald, 2013).
743 The discussion drew comparisons between the explicit focus on the teaching and
744 learning of domain-specific language against a growing perception that the culture
745 and practice of EAP are moving away from a focus on language towards generic skills,
746 and the implications that this shift in focus might have for teachers and students:

747 Knowledge User 6: I think one major, major, major issue with *EAP* is that it has become
748 so un-*language* focused. It's moved so far away from *teaching language*. And, *students*, of
749 course, can't understand this because that's what they think they're paying for. They think
750 we're there to *teach* them the *English*. I think I'm there to *teach* them the *English* but the
751 powers that be think that we're there to *teach* them *EAP*.

752 Researcher 1: I mean we didn't do any, there was no *time* in the timetables for *language*,
753 right?

754 Knowledge User 6: No, for *language*, nothing. It's all just skills.

755 Knowledge User 5: I couldn't believe it when I started *teaching EAP*.

756 Knowledge User 6: Skills and process. And this is so deeply concerning when they don't
757 have the *language* to express their ideas.

758 Knowledge User 5: I think that's why when they started this redundancy thing, oh well, I
759 didn't fight it because I'm not *teaching language* in *EAP* and I enjoy *teaching language*.

760 (Focus-group discussion excerpt with Knowledge User 5 & Knowledge User 6, Café Nero,
761 Durham UK)

762 Corpora provide teachers and learners with access to linguistic data that show
763 how language is used across a variety of real-world communication contexts. There
764 have been many successful commercial language coursebook publications that are
765 informed by corpora. However, many more coursebook publications appear to fly in
766 the face of evidence-based approaches to materials writing for meeting the demands
767 of an English language education content industry that seems to be driven, first and
768 foremost, by market research rather than research into whether or not materials have
769 positively influenced teaching, learning and language acquisition. A meeting with
770 EAP teacher, Knowledge User 1, highlights some of the issues with EAP materi-
771 als writing with commercial publishers. Despite materials not always drawing on
772 evidence of how language actually works, they are still widely marketed for sales
773 distribution and consumption:

774 Knowledge User 1: What I saw with him [EAP materials writer with Oxford University
775 Press] was, with his presentation at IATEFL [International Association for Teaching English
776 as a Foreign Language] was, that it was no more or less like really saying that THESE
777 *materials* he is selling are THE exponents that we *need to teach students*. And it was still

778 very much along the lines of we *need to teach* them yet more fixed phrases. And I was like
 779 sitting there and thinking some yes, some no, but prove it. I can—Can you? And he was
 780 putting up his *examples*, and I had my tablet *open* using *FLAX*, and I was going that *example*
 781 of his *works*, and that *works*, that doesn't *work*, that *works*, that doesn't *work*. But he's just
 782 basing it on his own judgement. And I'm just sitting there testing. Just right in front of him,
 783 testing his *materials*.

784 Researcher 1: And, you would have thought that he would have tested his *examples* with
 785 a corpus-informed approach before presenting them at IATEFL let alone publishing them
 786 with OUP. You have to wonder where the quality control lies if at all.

787 [...]

788 Knowledge User 1: The vast majority of my colleagues at Queen Mary have been pretty
 789 *open-minded*, and they've been looking at *FLAX* and they can see that it's real *academic*
 790 *language* data. It's the authenticity of it.

791 Researcher 1: Yes, that always wins out, doesn't it?

792 Knowledge User 1: Of course, it does but first of all they *need* to know that these non-
 793 commercial data-driven systems exist and that's where the commercial publishers have the
 794 upper hand.

795 (Meeting excerpt with Knowledge User 1, Cutty Sark pub in Greenwich, London UK)

796 6.5 Conclusion

797 With initiatives in open access and the changes to copyright legislation that have
 798 brought about TDM limitations and exceptions, we have seen the greatest distance
 799 travelled with this design-based research, resulting in the co-creation of the following
 800 language learning collections that remix open access content for learning features
 801 of academic English: the largest English language collocations collections used by
 802 learners online (Wu et al., 2021), the full-text BAWE collections in collaboration with
 803 EAP teachers at Durham University, the EThOS PhD abstract corpora with partici-
 804 pating EAP practitioners from Queen Mary University of London, the legal English
 805 BLaRC collection by Dr Maria Jose Marín from the University of Murcia, and the
 806 Academic Collocations in English (ACE) corpora with the COncnecting REpositories
 807 (CORE) aggregation and Application Programming Interface (API) services at the
 808 UK Open University. There is a growing sense that knowledge organisations such as
 809 the British Library and the Oxford Text Archive and aggregation and API services
 810 such as CORE are interested in non-commercial educational re-use applications
 811 of open access content that are aligned with the Budapest Open Access Initiative.
 812 Indeed, by far the biggest impact of openness in the higher education sector has been
 813 with open access, showing the importance of knowledge organisations in promoting
 814 accessible and reusable research (Finch Group, 2012).

815 The research presented on remixing MOOC content with TDM approaches
 816 provides proof of concept for the importance of licensing MOOC content openly
 817 for much-needed data-driven support with domain-specific language in non-formal

818 education that has re-use value in formal EAP education (Fitzgerald et al., 2017).
819 This increased value from open language learning online is echoed by the upswing in
820 the enrolment in language MOOCs that have emerged during the pandemic (Martín-
821 Monje & Borthwick, 2021). However, findings from our research point to a current
822 problem with the scalability of developing derivative OERs from MOOC content,
823 with the example presented here of providing data-driven language support in the
824 MOOC context. This problem is apparent in current mainstream MOOC provision
825 where current business models do not anticipate a need for the open licensing of
826 course content, and where open educational practices are mostly limited to those
827 subject academics and learning technologists who were already open digital scholars
828 before engaging in MOOC and networked learning pedagogy. Rather, current MOOC
829 business models appear to focus on charging learners for increased access to learning
830 content. This phenomenon has been presented here as an issue that open education
831 policy makers, in collaboration with Creative Commons, are actively lobbying to
832 address. As a work-around solution for embedding the functions and open corpora
833 of FLAX directly into a MOOC platform interface, research is currently being carried
834 out by Dr Jemma König at the University of Waikato with the development of F-
835 Lingo, a Chrome extension. F-Lingo works on top of the FutureLearn platform to
836 support content-based learning of domain-specific terminology and concepts for
837 academic and professional English. Nonetheless, this work with F-Lingo would still
838 require higher education institutions to allow the traversing and re-use of All Rights
839 Reserved course content for the R&D of automated language learning support in the
840 MOOC context (Fitzgerald et al., 2019; König et al., 2022).

841 The observed absence of data-driven approaches to support EAP provision at two
842 UK universities, and the apparent shift away from language teaching, as noted in
843 focus-group discussions with teachers and managers, give pause for understanding
844 current practices with EAP materials development for classroom and online learning
845 in a time of increased uptake of generic EAP course books from commercial
846 publishers. By drawing attention to the underlying business models and cultural
847 practices that higher education institutions and organisations adopt, we also arrive at
848 a closer understanding of the values placed on research, or lack thereof, with online
849 and classroom materials development and teaching in the field of EAP.

850 The new paradigm for open data-driven language learning systems design
851 presented through this research has also argued for greater access to and re-use of the
852 artefacts of the academy and professional domains such as law, for example, that are
853 taught and studied at higher education institutions. In this chapter, we have demon-
854 strated the perceived value that corpus linguistics researchers and knowledge users
855 working within EAP place on pedagogic, professional and research texts that can be
856 mined for aspects of domain-specific language with data-driven learning systems. In
857 addition to the value placed on open educational practices that can be fostered to re-
858 use, remix and redistribute EAP resources for uptake across formal and non-formal
859 higher education in increasingly uncertain times.

860 **Acknowledgements** We would like to thank the many contributors and collaborators of this
861 ongoing R&D project over the years, including the Fonds de recherche du Québec—Société et

862 culture (FRQSC), the OER Research Hub and the Global OER Graduate Network based at the UK
 863 Open University, and the International Research Foundation (TIRF) for English Language Education
 864 for funding this research collaboration between the FLAX and F-Lingo projects at the Department
 865 of Computer Science at the University of Waikato in Aotearoa New Zealand, the Department of
 866 Education at Concordia University in Montréal Canada, the Departamento de Filología Inglesa at
 867 Universidad de Murcia in Spain, the School of Languages, Linguistics and Film at Queen Mary
 868 University of London and the Durham Centre for Academic Development at Durham University in
 869 the UK.

870 Appendix

871 Open Collections in FLAX: Content and Collaborators

Learning collocations system in FLAX (2009–2022)

Content	<ul style="list-style-type: none"> • Wikipedia corpus of contemporary English derived from three million Wikipedia articles comprising three billion words (Wu & Witten, 2016; Wu et al., 2021) • British National Corpus (BNC) of 100 million words (BNC Consortium, 2007) • British Academic Written English (BAWE) corpus of 2500 pieces of assessed university student writing from across the disciplines • Academic Collocations in English (ACE) corpora of harvested open access content and metadata from 135 million articles residing in open journals and open repositories
Knowledge organisations	Wikimedia Foundation (Wikipedia corpus); Oxford Text Archive and the UK Higher Education Academy OER International Programme with the University of Oxford (BNC and BAWE corpora); CORE (Connecting REpositories) ¹⁰ team, UK Open University (ACE corpora)
Researchers	FLAX team
Knowledge users	Waikato University computer science students; Durham University EAP teachers and students; University of Oxford OER International stakeholders

British Academic Written English (BAWE) collections in FLAX (2012)

Content	Full texts of the BAWE corpus divided into four sub-collections: Arts & Humanities, Social Sciences, Life Sciences, Physical Sciences
Knowledge organisations	The Oxford Text Archive; UK Higher Education Academy
Researchers	FLAX team
Knowledge users	Durham University EAP teachers and students; University of Oxford OER International stakeholders

British Law Report Corpus (BLaRC) in FLAX (2014)

(continued)

¹⁰ <https://core.ac.uk/about#mission>.

(continued)

Content	8.85 million-word corpus of full-text judicial hearings derived from free legal sources at the British and Irish Legal Information Institute (BAILII) ¹¹ aggregation website
Knowledge organisations	BAILII
Researchers	Universidad Murcia; FLAX team
Knowledge users	Law MOOC learners

MOOC/micro-networked course collections in FLAX (2014–2016)

Content	MOOC / Micro-Networked Course lecture transcripts and videos (streamed via YouTube or Vimeo) and case law that reside in the public domain
Knowledge organisations	MOOC host institutions (Harvard University; University of London; Columbia University) with edX and Coursera MOOC providers
Researchers	FLAX team; Universidad Murcia
Knowledge users	MOOC learners and MOOC subject matter experts; legal English translation studies teachers, and students at the University of Murcia

PhD micro-abstract corpora with FLAX mobile activities (2014–2015)

Content	Domain-specific micro abstract corpora in the areas of Law, Water Politics and Tourism Studies. Developed in collaboration with EAP teachers at Queen Mary University of London for use on summer EAP pre-sessional courses. Developed with web-based and mobile language learning activities using the suite of mobile applications for Android from FLAX
Knowledge organisations	British Library Labs ¹² and ETHOS ¹³ at the British Library
Researchers	FLAX team
Knowledge users	EAP teachers and learners at Queen Mary University of London

PhD abstract corpora in FLAX (2015–2016)

Content	9.8 million-word corpus derived from the metadata, including the abstracts, of over 500,000 PhD theses awarded by UK universities and managed by the Electronic Thesis Online Service (ETHOS) at the British Library
Knowledge organisations	British Library Labs and ETHOS at the British Library
Researchers	FLAX and F-Lingo teams
Knowledge users	EAP teachers and managers at Queen Mary University of London; Current research with MOOC learners via F-Lingo Chrome extension and FutureLearn platform

Academic Collocations in English (ACE) collections in FLAX (2018–2022)

Content	Harvested open access content from open journals and open repositories divided into four sub-collections: Arts & Humanities, Social Sciences, Life Sciences, Physical Sciences
---------	--

(continued)

¹¹ <http://ials.sas.ac.uk/digital/bailii>.¹² <https://www.bl.uk/projects/british-library-labs>.¹³ <http://ethos.bl.uk/Home.do>.

(continued)

Knowledge organisations	CORE (Connecting REpositories) team, UK Open University
Researchers	FLAX and F-Lingo teams
Knowledge users	<ul style="list-style-type: none"> • User query data analysis research with the FLAX LC system learners worldwide • Research with MOOC learners via F-Lingo Chrome extension and FutureLearn platform

References

- Amiel, T., & Reeves, T. C. (2008). Design-based research and educational technology: Rethinking technology and the research agenda. *Educational Technology & Society*, 11(4), 29–40.
- Anderson, T., & Shattuck, J. (2012). Design-based research: A decade of progress in education research? *Educational Researcher*, 41(1), 16–25.
- Angus, D., Rintel, S., & Wiles, J. (2013). Making sense of big text: A visual-first approach for analyzing text data using Leximancer and Discursis. *International Journal of Social Research Methodology*, 16(3), 261–267.
- Anthony, L. (2014, July). A view to the future in corpus tools development. In *11th Teaching and Language Corpora Conference (TALC 11) Keynote Address*. Lancaster University, UK.
- Anthony, L. (2018). *Introducing English for specific purposes*. Routledge.
- Atenas, J., Havemann, L., & Priego, E. (2015). Open data as open educational resources: Towards transversal skills and global citizenship. *Open Praxis*, 7(4), 377–389.
- Barab, S., & Squire, L. (2004). Design-based research: Putting a stake in the ground. *The Journal of the Learning Sciences*, 13(1), 1–14.
- Belcher, D. (2010). What ESP is and can be: An introduction. In D. Belcher (Ed.), *English for specific purposes in theory and practice* (pp. 1–20). University of Michigan Press.
- Blei, D. (2012a). Probabilistic topic models. *Communications of the ACM*, 55, 77–84.
- Blei, D. (2012b). Topic modelling and digital humanities. *Journal of Digital Humanities*, 2, 8–11.
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3, 993–1022.
- Boulton, A., & Thomas, J. (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* (pp. 7–34). Masaryk University Press.
- Boulton, A., & Cobb, T. (2017). Corpus use in language learning: A meta-analysis. *Language Learning*, 67(2), 348–393.
- Boulton, A., & Pérez-Paredes, P. (2014). ReCALL special issue: Researching uses of corpora for language teaching and learning editorial researching uses of corpora for language teaching and learning. *ReCALL*, 26(2), 121–127.
- Boumans, J. W., & Trilling, D. (2016). Taking stock of the toolkit: An overview of relevant automated content analysis approaches and techniques for digital journalism scholars. *Digital Journalism*, 4(1), 8–23.
- British Library. (n.d.). EThOS Toolkit | Re-use by external services: EThOS as a data provider: Metadata. Retrieved from <http://ethostoolkit.cranfield.ac.uk/tiki-index.php?page=Re-use+by+external+services>
- Budapest Open Access Initiative (BOAI). (2002). BOAI declaration. Retrieved from <http://www.budapestopenaccessinitiative.org/read>
- Burns, A. (2009). *Doing action research in English language teaching*. Routledge.

- 910 Chang, J.-Y. (2014). The use of general and specialized corpora as reference sources for academic
911 English writing: A case study. *ReCALL*, 26(2), 243–259.
- 912 Charles, M. (2012). Proper vocabulary and juicy collocations: EAP students evaluate do-it-yourself
913 corpus-building. *English for Specific Purposes*, 31(2), 93–102.
- 914 Charles, M. (2015). Same task, different corpus: The role of personal corpora in EAP classes.
915 In A. Lenko-Szymanska & A. Boulton (Eds.), *Multiple affordances of language corpora for*
916 *data-driven learning* (pp. 131–153). John Benjamins.
- 917 Charles, M., & Pecorari, D. (2016). *Introducing English for academic purposes*. Routledge.
- 918 Cobb, T., & Boulton, A. (2015). Classroom applications of corpus analysis. In D. Biber & R. Reppen
919 (Eds.), *Cambridge handbook of English corpus linguistics* (pp. 478–497). Cambridge University
920 Press.
- 921 Colpaert, J. (2016). Big content in an educational engineering approach. *Journal of Technology*
922 *and Chinese Language Teaching*, 7(1), 1–14. Retrieved from [http://tclt.us/journal/2016v7n1/col](http://tclt.us/journal/2016v7n1/colpaert.pdf)
923 [paert.pdf](http://tclt.us/journal/2016v7n1/colpaert.pdf)
- 924 Colpaert, J. (2004). Transdisciplinarity. *Computer Assisted Language Learning*, 17(5), 459–472.
- 925 Colpaert, J. (2018). Transdisciplinarity revisited. *Computer Assisted Language Learning*, 31(5–6),
926 483–489.
- 927 COnnectedREpositories (CORE): Aggregating the world's open access papers. (n.d.). Retrieved
928 from <https://core.ac.uk/>
- 929 Ding, A., & Bruce, I. (2017). *The English for academic purposes practitioner: Operating on the*
930 *edge of academia*. Palgrave Macmillan.
- 931 ETIC. (1975). *English for academic study: Problems and perspectives*. British Council.
- 932 Finch Group. (2012). Accessibility, sustainability, excellence: How to expand access to research
933 publications. *Report of the Working Group on Expanding Access to Published Research Findings*.
934 Retrieved from <http://www.researchinfonet.org/publish/finch/>
- 935 Fitzgerald, A. (2013). *Openness in English for academic purposes*. Open Educational Resources
936 Case Study: Pedagogical development from OER practice. Commissioned by the Higher Educa-
937 tion Academy (HEA) and the Joint Information Systems Committee (JISC), United Kingdom,
938 20 pages.
- 939 Fitzgerald, A., Wu, S., & Barge, M. (2014). Investigating an open methodology for designing
940 domain-specific language collections. In S. Jager, L. Bradley, E. J. Meima & S. Thouésny (Eds.),
941 CALL design: Principles and practice. In *Proceedings of the 2014 EUROCALL Conference*
942 (pp. 88–95). Groningen, The Netherlands. Dublin: Research-publishing.net. [https://doi.org/10.](https://doi.org/10.14705/rpnet.2014.000200)
943 [14705/rpnet.2014.000200](https://doi.org/10.14705/rpnet.2014.000200).
- 944 Fitzgerald, A., Wu, S. & Marin, M.J. (2015). FLAX—Flexible and open corpus-based language
945 collections development. In K. Borthwick, E. Corradini & A. Dickens (Eds.), *10 years of the*
946 *Languages, Linguistics & Area Studies (LLAS) eLearning symposium: case studies in good*
947 *practice* (pp. 215–227). Research-publishing.net. <https://doi.org/10.14705/rpnet.2015.000281>
- 948 Fitzgerald, A., Marín, M. J., Wu, S., & Witten, I. H. (2017). Evaluating the efficacy of the digital
949 commons for scaling data-driven learning. In M. Carrier, R. Damerow & K. Bailey (Eds.), *Digital*
950 *language learning and teaching: Research, theory and practice* (pp. 38–51). Global Research on
951 Teaching and Learning English Series. Routledge, Taylor & Francis.
- 952 Fitzgerald, A., König, J., & Witten, I. H. (2019). F-Lingo: Integrating lexical feature identifica-
953 tion into MOOC platforms for learning professional and academic English. In R. Meir, J.
954 Sluss, E. Tovar & M. Castro (Eds.), *Proceedings of the 6th Conference on Learning with*
955 *MOOCs: Enhancing Workforce Diversity and Inclusion* (pp. 101–104). Institute of Electrical
956 and Electronics Engineers (IEEE) Education Society.
- 957 Fitzgerald, A. (2019). *A new paradigm for open data-driven language learning systems design in*
958 *higher education*. Unpublished doctoral thesis. Concordia University, Canada.
- 959 Flowerdew, J., & Peacock, M. (2001). Issues in EAP: A preliminary perspective. In J. Flow-
960 erdew & M. Peacock (Eds.), *Research perspectives on English for academic purposes* (pp. 8–24).
961 Cambridge University Press.

- 962 Gillett, A. (2018, April 13). Is EAP ESP? [Blog post] *Uefap*. Retrieved from [http://www.uefap.net/](http://www.uefap.net/blog/?p=933)
 963 [blog/?p=933](http://www.uefap.net/blog/?p=933)
- 964 Hadley, G. (2015). *English for academic purposes in neoliberal universities: A critical grounded*
 965 *theory* (Vol. 22). Springer International Publishing.
- 966 Hakkarainen, P. (2009). Designing and implementing a PBL course on educational digital video
 967 production: Lessons learned from a design-based research. *Educational Technology, Research*
 968 *and Development*, 57(2), 211–228. <https://doi.org/10.1007/s11423-007-9039-4>
- 969 Hargreaves, I. (2011). *Digital opportunity—A review of intellectual property and growth*. HM
 970 Government. Retrieved from [https://assets.publishing.service.gov.uk/government/uploads/sys](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/32563/ipreview-finalreport.pdf)
 971 [tem/uploads/attachment_data/file/32563/ipreview-finalreport.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/32563/ipreview-finalreport.pdf)
- 972 Herrington, J., McKenney, S., Reeves, C. T., & Oliver, R. (2007). Design-based research and doctoral
 973 students: Guidelines for preparing a dissertation proposal. Retrieved from [http://ro.ecu.edu.au/](http://ro.ecu.edu.au/cgi/viewcontent.cgi?article=2611&context=ecuworks)
 974 [cgi/viewcontent.cgi?article=2611&context=ecuworks](http://ro.ecu.edu.au/cgi/viewcontent.cgi?article=2611&context=ecuworks)
- 975 Howatt, A. P. R. (2004). *A history of English language teaching* (2nd ed.). Oxford University Press.
- 976 Hunston, S. (2002). *Corpora in applied linguistics*. Cambridge University Press.
- 977 Hyland, K. (2002). Specificity revisited: How far should we go now? *English for Specific Purposes*,
 978 21, 385–395.
- 979 IT Services, University of Oxford, 13 Banbury Road. (n.d.). [Oxford Text Archive] [BAWE Terms
 980 and Conditions Text]. Retrieved from <http://ota.ox.ac.uk/scripts/download.php?otaid=2539>
- 981 Johns, T. (1991b). From printout to handout: Grammar and vocabulary teaching in the context of
 982 data-driven learning. In T. Johns & P. King (Eds.), *Classroom concordancing*. *English Language*
 983 *Research Journal*, 4, 27–45.
- 984 Johns, T. (1991a). Should you be persuaded: two examples of data-driven learning. In T. Johns &
 985 P. King (Eds.), *Classroom concordancing*. *English Language Research Journal*, 4, 1–16.
- 986 Johns, T. (2002). Data-driven learning: The perpetual challenge. In B. Kettemann & G. Marko
 987 (Eds.), *Teaching and learning by doing corpus analysis. Proceedings of the Fourth International*
 988 *Conference on Teaching and Language Corpora, Graz 19–24 July 2000* (pp. 107–117). Rodopi.
- 989 König, J., Wu, S., Fitzgerald, A., Franken, M., & Witten, I. H. (2022). F-Lingo: Leveraging smart
 990 CALL for massive open online courses. In J. Colpaert & G. Stockwell (Eds.), *Smart CALL*,
 991 Waseda University.
- 992 Levin, B. (2011). Mobilising research knowledge in education. *London Review of Education*, 9(1),
 993 15–26.
- 994 Marín, M. J., & Rea, C. (2014). Assessing four automatic term recognition methods: Are they
 995 domain-dependent? *English for Specific Purposes World*, 42(15), 1–27.
- 996 Martín-Monje, E., & Borthwick, K. (2021). Researching massive open online courses for language
 997 teaching and learning. *ReCALL*, 33(2), 107–110.
- 998 McEnery, T., Xiao, R., & Tono, Y. (2006). *Corpus-based language studies: An advanced resource*
 999 *book*. Routledge.
- 1000 McKenney, S., & Reeves, T. (2012). *Conducting educational design research*. Routledge.
- 1001 Nunez-Mir, G. C., Iannone, B. V., Pijanowski, B. C., Kong, N., & Fei, S. (2016). Automated content
 1002 analysis: Addressing the big literature challenge in ecology and evolution. *Methods in Ecology*
 1003 *and Evolution*, 7(11), 1262–1272.
- 1004 Papadimitriou, C. H., Tamaki, H., Raghavan, P., & Vempala, S. (1998). Latent semantic indexing:
 1005 A probabilistic analysis. In *Proceedings of the Seventeenth ACM SIGACT-SIGMOD-SIGART*
 1006 *Symposium on Principles of Database Systems* (pp. 159–168). Seattle, Washington.
- 1007 Pérez-Paredes, P., Ordoñana Guillamón, C., & Aguado Jiménez, P. (2018). Language teachers’
 1008 perceptions on the use of OER language processing technologies in MALL. *Computer Assisted*
 1009 *Language Learning*, 31(5–6), 522–545.
- 1010 Reason, P., & Bradbury, H. (2007). *Handbook of Action Research* (2nd ed.). Sage.
- 1011 Sinclair, J. M. (2004). *Trust the text: Language, corpus and discourse*. Routledge.
- 1012 Smith, A. E. (2000b). Machine mapping of document collections: The Leximancer system. In
 1013 *Proceedings of the Fifth Australasian Document Computing Symposium*. DSTC.

- 1014 Smith, A. E. (2000a). Machine learning of well-defined thesaurus concepts. In A.-H. Tan & P. S.
 1015 Yu (Eds.), *Proceedings of the International Workshop on Text and Web Mining (PRICAI 2000a)*
 1016 (pp. 72–79). PRICAI.
- 1017 Smith, A. E. (2003). Automatic extraction of semantic networks from text using Leximancer. In
 1018 *HLT-NAACL 2003 Human Language Technology Conference of the North American Chapter*
 1019 *of the Association for Computational Linguistics: Companion Volume* (pp. Demo23–Demo24).
 1020 ACL.
- 1021 Smith, A. E., & Humphreys, M. S. (2006). Evaluation of unsupervised semantic mapping of natural
 1022 language with Leximancer concept mapping. *Behavior Research Methods*, 38(2), 262–279.
- 1023 Strevens, P. (1988). ESP after twenty years: A reappraisal. In M. Tickoo (Ed.), *ESP: State of the*
 1024 *art* (pp. 1–13). SEAMEO Regional Language Centre.
- 1025 Swales, J. M. (1985). *Episodes in ESP*. Pergamon Press.
- 1026 Thomas, J. (2017). *Discovering English with sketch engine* (2nd ed.). Versatile.
- 1027 Turner, J. (2004). Language as academic purpose. *Journal of English for Academic Purposes*, 3,
 1028 95–109.
- 1029 Vyatkina, N. (2016). Data-driven learning of collocations: Learning performance, proficiency, and
 1030 perceptions. *Language Learning & Technology*, 20(3), 159–179.
- 1031 Walker, D. (2006). Towards productive design studies. In J. van den Akker, K. Gravemeijer, S.
 1032 McKenney, & N. Nieveen (Eds.), *Educational design research* (pp. 9–18). Routledge.
- 1033 Weller, M. (2011). *The digital scholar: How technology is changing scholarly practice*. Bloomsbury
 1034 Academic.
- 1035 Widdowson, H. G. (1983). *Learning purpose and language use*. Oxford University Press.
- 1036 Wu, S., Fitzgerald, A., Witten, I. H., & Yu, A. (2018). Automatically augmenting academic text for
 1037 language learning: PhD abstract corpora with the British Library. In B. Zou & M. Thomas (Eds.),
 1038 *Integrating technology into contemporary language learning and teaching* (pp. 512–537), IGI
 1039 Global.
- 1040 Wu, S., Fitzgerald, A., Yu, A., & Chen, Z. (2021). What are language learners looking for in a
 1041 collocation consultation system? Identifying collocation look-up patterns with user query data.
 1042 *ReCALL*, 33(3), 229–247.
- 1043 Wu, S., & Witten, I. H. (2016). Transcending concordance: Augmenting academic text for L2
 1044 writing. *International Journal of Computer-Assisted Language Learning and Teaching*, 6(2),
 1045 1–18.

1046 **Alannah Fitzgerald** is a postdoctoral research fellow with the Computer Science Department at
 1047 the University of Waikato in Aotearoa, New Zealand, and an honorary fellow with the School
 1048 of Education at Durham University in the UK. Alannah is responsible for designing open educa-
 1049 tion applications with the F-Lingo and FLAX language projects. With Dr. Wu, Professor Witten
 1050 and Chris Mansfield, she was awarded a prize in the British Library Labs Competition Teaching
 1051 and Learning category for re-using digital collections in language education with funding from
 1052 the Andrew W. Mellon Foundation. Her research interests include open educational resources and
 1053 practices for designing and developing digital domain-specific language collections (corpora) and
 1054 for devising and delivering online English language learning interventions that can be scaled and
 1055 assessed across both formal (classroom-based) and non-formal (MOOC space) higher education
 1056 contexts.

1057 **Shaoqun Wu** is Senior Lecturer in Computer Science at the University of Waikato in Aotearoa,
 1058 New Zealand and is the main developer of the FLAX language project. With Dr. Fitzgerald and
 1059 Professor Witten, Dr. Wu was awarded first prize in the LinkedUp Vici Challenge for mature
 1060 open data-driven applications for education by Open Knowledge International with funding from
 1061 the European Commission. Her research interests include computer-assisted language learning,
 1062 mobile language learning, supporting language learning in MOOCs, digital libraries, natural
 1063 language processing and computer science education.

1064 **Jemma König** is Postdoctoral Fellow in the Department of Computer Science at the University
1065 of Waikato in Aotearoa, New Zealand and is responsible for developing the F-Lingo Chrome
1066 extension for FutureLearn MOOCs. Jemma's PhD research explored a computational approach
1067 to vocabulary testing, language tools and text enrichment. More specifically, focusing on corpus
1068 analysis, pseudoword generation, automated vocabulary testing and tracking learners' interaction
1069 with online written language. With Dr. Fitzgerald and Professor Witten, Dr. König was awarded
1070 best paper for her work with F-Lingo at the Learning with MOOCs conference in Milwaukee,
1071 USA, by the Institute of Electrical and Electronics Engineers (IEEE) Education Society in 2019.

1072 **Steven Shaw** is a professor in the Department of Education, Concordia University, in Montreal.
1073 His research and professional work focus on the design, development and implementation of tech-
1074 nology to support learning and knowledge sharing, particularly at the enterprise scale in large
1075 public and private sector organisations. He co-founded the corporation that developed the first
1076 "learning content management system". He served as the CLO of Eedo Knowledgeware, which for
1077 over a decade furnished the leading-edge technology for learning content management, employed
1078 by Fortune 500 organisations such as Xerox, Dell, Eli Lilly, Boeing and the largest public sector
1079 organisations in the US and UK, including US Treasury and Department of Energy and Foreign
1080 and Commonwealth Office and Department of Work and Pensions in the UK. His areas of
1081 expertise include software development, systems implementation, content management, taxonomy
1082 development and the design and evaluation of training programs and curricula in professional
1083 education.

1084 **Ian H. Witten** is Emeritus Professor of Computer Science at the University of Waikato in
1085 Aotearoa, New Zealand, with a research career that spans over 40 years. His best-known publi-
1086 cation is the book, *Data mining: Practical machine learning tools and techniques*, now in its
1087 fourth edition (2016). Professor Witten is also well known for his award-winning open-source
1088 software, sharing his advances with thousands of students, teachers and users around the world.
1089 These include Greenstone, a digital library platform on which the FLAX system operates. Another
1090 successful open-source software is Weka (Waikato Environment for Knowledge Analysis), a data-
1091 mining tool. Weka is probably the world's most widely used machine learning workbench, and in
1092 2017 Professor Witten led three popular Massive Open Online Courses (MOOCs) with Future-
1093 Learn *Data mining with Weka*, *More data mining*, and *Advanced data mining*. In 2017, Professor
1094 Witten was awarded an Honorary PhD (Doctor of the University) from the Open University in
1095 the United Kingdom for his lifetime contribution to furthering the advancement of research and
1096 education.

Author Queries

Chapter 6

Query Refs.	Details Required	Author's response
AQ1	Please note that the footnote has been set in the following sentence "It identifies the knowledge organisations,...", as footnotes are not allowed in abstracts.	
AQ2	References "Burkhardt (2006), Weber (1990), Fisher (2014), BNC Consortium (2007)" are cited in the text but not provided in the reference list. Please provide the respective references in the list or delete these citations.	