

## A DATA-DRIVEN LEARNING EXPERIMENT IN THE LEGAL ENGLISH CLASSROOM USING THE FLAX PLATFORM

**MARÍA JOSÉ MARÍN**

Universidad de Murcia  
Mariajose.marin1@um.es

**MARÍA ÁNGELES ORTS LLOPIS**

Universidad de Murcia  
mageorts@um.es

**ALANNAH FITZGERALD**

Concordia University, Montreal  
alannahfitzgerald@gmail.com

---

37

### Abstract

This research presents a data-driven experiment in the legal English field where the FLAX, an open-source self-learning online platform, is assessed as regards its efficacy in aiding a group of legal English non-native undergraduates (divided into an experimental and a control group) to use legal terminology more consistently, amongst other language items. The experimental group were instructed to only resort to the FLAX and to exploit all the functionalities offered by it. Conversely, the control group could access any information source at hand except for the learning platform for the completion of the same task. Two learner corpora were gathered and analysed on a lexical and pragmatic level for the evaluation of term usage and distribution, lexical diversity, lexical fundamentality and the use of discourse markers. The results display a tendency on the part of the experimental group towards a more consistent usage of legal terminology, which also appears to be better distributed than the terms in the non-FLAX corpus. In contrast and on average, the lexicon in the FLAX-based corpus tends to be slightly more basic. Concerning the use of MD markers, the experimental group appears to use, though marginally, a greater number of evidentials, endophoric and interactional markers.

**Keywords:** legal English, data-driven learning (DDL), corpus linguistics, learner corpora, open access.

## Resumen

En este artículo se presenta un experimento basado en corpus para la enseñanza del inglés jurídico donde se evalúa la plataforma FLAX, un sistema online de aprendizaje de lenguas, como apoyo a la enseñanza de esta variedad del inglés. Los informantes fueron divididos en un grupo experimental y otro de control. Al grupo experimental se le pidió que utilizara únicamente FLAX para la realización de la tarea haciendo uso de todas las opciones que facilita dicha plataforma. Por el contrario, el grupo de control podría utilizar cualquier fuente de información para la realización del trabajo a excepción de FLAX. Se compilaron dos corpus con el material elaborado por los informantes y se analizaron a nivel léxico y pragmático para la evaluación del uso y la distribución de la terminología especializada, la diversidad léxica y el uso de los marcadores del discurso. Los resultados muestran una tendencia por parte del grupo experimental hacia un uso más consistente de la terminología jurídica, que además parece estar mejor distribuida que lo está en el corpus del grupo de control. En lo que respecta al uso de los marcadores del discurso, el grupo experimental emplea un mayor número de marcadores endofóricos, interaccionales y evidenciales.

**Palabras clave:** inglés jurídico, data-driven learning (DDL), lingüística del corpus, learner corpora, open access.

## 1. Introduction

The use of language corpora in language instruction has been explored profusely, as illustrated by authors like Boulton (2010a), since they can contribute, not only to the provision of authentic language samples which enable learners and instructors to approach language learning from a different perspective, but also to the learning process itself. As Johns (1986; 1991; 1997) –who coins the term *data-driven learning* (DDL henceforth) – points out, through the direct observation of corpus samples, students can infer the rules of language and “develop strategies for discovery –strategies through which he or she can learn how to learn–” (Johns 1991: 1). In other words, they can become “language detectives” (Johns 1997: 101).

There exists a large number of teaching resources focused on general English basically due to the number of potential users of these teaching materials and the economic benefits this might generate. However, and precisely due to that fact, the more specialised the need, the fewer materials we find, as Boulton (2012) acknowledges. As regards corpus-based materials specifically, some scholars (McEnery and Wilson 1996; Boulton 2010a) consider that they address the

students' needs better than other traditional materials like coursebooks, "including quantitative accounts of vocabulary and usage which address the specific needs of students in a particular domain more directly than those taken from more general language corpora" (McEnery and Wilson 1996: 121). In Boulton's words (2012: 262), they can provide "a framework to highlight the highly conventionalised language used in specialist disciplines, especially where the focus is on a specific genre or text type".

As a consequence of this tendency, there is a plethora of studies aimed at testing the efficiency and advantages/disadvantages of corpus-based language instruction within the general and specific fields, yet, unlike other ESAP (English for Specific and Academic Purposes) varieties, legal English has not been sufficiently explored or tested in this respect (Boulton 2010b; Marín 2014b; Marín and Fernández Toledo 2015). This was one of the major reasons which motivated the present research, which aims at determining the efficiency and influence of corpus-based materials on the usage of legal English terminology, at a lexical level, and the expression of engagement and stance through the use of metadiscourse markers, at a pragmatic level.

To that end, two learner corpora were gathered, which comprised the essays written by 105 undergraduate students (divided into an experimental and a control group) as part of the final assessment of their legal English translation course. The essays presented the structure of research articles where the informants had to critically review the literature related to each topic of their choice. Among the topics they had to write about were contract law, international law, common vs. civil law, the sources of law, the principle of binding precedent, legal genres, criminal law: major offences, or probate law.

For the completion of this task, the experimental group was only allowed to consult and exploit the different functionalities offered by the FLAX,<sup>1</sup> a corpus-based open-source language platform, while the control group could refer to any information source at hand. The essays were then processed applying corpus linguistics techniques, which allowed us to quantify term usage and distribution, lexical diversity and fundamentality and also to reveal interpersonality traits based on an analysis of metadiscourse markers.

Two research questions were thus formulated:

RQ1: Would this corpus-based platform positively influence the usage of specialised legal terminology by learners?

RQ2: Can corpus-based materials also influence the usage of metadiscourse markers?

## 2. Literature review

The potential benefits of the use of language corpora in second language teaching and learning have been discussed by scholars such as Johns (1986, 1991), Sinclair (1991, 2003), McEnery and Wilson (1996), McEnery and Xiao (2011), Hunston (2007) or Boulton (2011), to name but a few, who, amongst other advantages, highlight their capacity to present learners with authentic materials and to offer plenty of genuine examples of a particular linguistic item in various contexts, thus facilitating its understanding through such contexts. Not only can corpora assist understanding through contextualisation and offer samples of the language in authentic settings, but they can also contribute to learners' motivation, as initially put forward by Johns (1986, 1991) and later by Boulton (2011), who affirms that they are capable of "empowering learners to explore language corpora and come to their own conclusions" (2011: 563).

Nevertheless, one of the main criticisms levelled at DDL methods, according to some authors who follow the chomskian trend (less numerous than those who support their usage), is precisely related to the context of corpus language samples, which appears to be insufficient in their view. As Flowerdew (2009: 406) puts it, corpus samples, if selected at random and analysed in the SL classroom, are "truncated concordance lines [which] are examined atomistically". This is precisely why Hunston (2007) recommends that such samples should be filtered, selected and adapted to the students' levels and needs.

There have been many DDL experiments aimed at testing the efficiency of corpora in supporting second language acquisition processes in specialised settings like translation (Aston 1997), technical engineering (Todd 2001), economics (Hadley 2002), computing (Clerehan et al. 2003), tourism (Curado Fuentes 2004) or architecture (Boulton 2010a), to name but a few. In Boulton (2010b) we find a comprehensive review of over a hundred different empirical evaluations of DDL carried out in the last two decades. Yet, the legal English field remains underexplored as only two of these experiments are dedicated to this ESAP branch (Fan and Xun-Feng 2002; Hafner and Candlin 2007). The scenario is similar in EAP (English for Academic Purposes) writing, as Ädel (2010) points out, since there is a very limited number of studies implementing DDL methods in specialised or academic writing instruction. The direct approach, which, following Hunston and Römer (in Ädel 2010), consists in giving the students "hands-on access to corpora" in the SL classroom, appears to be the most controversial and also least explored DDL method, which poses a challenge for researchers working in the field. The present study falls within this category.

The research questions posed in the introduction to this study present two major foci, firstly, to measure the influence of resorting to a corpus-based learning platform on

the use of legal terminology by ESAP learners and, secondly, to try to access the pragmatic level of two learner corpora through the analysis of meta-discourse markers. Regarding the usefulness of meta-discourse markers (MD markers henceforth) in writing, Hyland (2005: 3) suggests that “the writer is not simply presenting information about the suggested route by just listing changes of direction, but taking the trouble to see the walk from the reader’s perspective”. Metadiscourse is, according to Hyland, “the means by which propositional content is made coherent, intelligible and persuasive” to receivers of texts (Hyland 2005: 39). MD markers could thus be regarded as tools for the expression of interpersonality, a concept that relates to Bakhtin’s/Voloshinov’s now widely influential notions of dialogism and heteroglossia. Interpersonality is a somewhat fuzzy concept that has been approached from different viewpoints such as the theories of appraisal, stance, evaluation and engagement (Biber and Finnegan 1989; Martin and White 2005; Sancho Guinda and Hyland 2012, among many others). Generally, the concept has been taken up and used by researchers to trace patterns of interaction and to discuss different aspects of language in use: the greater the abundance of markers, the clearer, the more legible and engaging the text is supposed to be.

The taxonomy for analysis deployed in this article will be based upon Hyland’s conception of MD (2005), the incidence of these markers in our texts being scrutinised in order to ascertain the level of proximity between interactants. Hyland organises metadiscourse markers by distinguishing between *interactive or textual devices* (those which organise information in an intelligible and persuasive way for the audience) and *interactional devices* (those that allow writers to articulate linguistically their attitudes and perspectives toward the propositional content of the text). In other words, through the use of textual markers, writers would be able to present the propositional content and their ideas both coherently and intelligibly to the readers, while interactive markers would, in turn, build an interaction between the reader and writer and create rapport and reader-friendliness in the text (Hyland and Tse 2004). The taxonomy of interactive or textual signals used by Hyland (2005; Hyland and Tse 2004) divides MD markers into transitions (conjunctions and conjunctives that help the readers determine the logical relationships between propositions), endophorics (referring to other parts of the text in order to make additional information available), frames (used to sequence parts of the text), glosses (supplying additional information by rephrasing, illustrating or explaining) and evidentials (helping to establish authorial command of the subject). According to Dafouz (2008), textual MD markers engage the reader on a level that relates more to formal grammar and are generally realised in the form of conjuncts and adverbials. The incidence of these markers in both our learner corpora will be quantified so as to measure how interpersonality is expressed in both text collections by resorting to them.

The textual function is intrinsic to language and exists to construe both propositional and interpersonal aspects into a linear and coherent whole. In comparison, interactional markers –hedges (indicators of the writer’s decision to recognise other voices), boosters (expressing authorial certainty), attitude markers (indicating the authorial opinion or assessment), engagement markers (drawing addressees into the discourse) and references to self (making authorial presence explicit in the text)– relate more to the socio-affective level where audience engagement from that perspective is prioritised in discourse (Heng and Tan 2010). The incidence of these markers in our texts will be scrutinised in order to ascertain the level of proximity between interactants, since, according to Mao (1993: 270), metadiscourse is not merely a stylistic device, but has a rhetorical role very much in line with the purpose that the text wishes to accomplish.

### **3. Methodology**

#### **3.1. The FLAX: an Open-Source Online Language Learning Platform**

42

The FLAX could be described as an open-source self-learning platform which mines salient linguistic features from augmented full-text corpora and displays them in interfaces designed to support learners with domain-specific language learning materials. Unlike traditional concordancers,<sup>2</sup> the FLAX project has developed interfaces for non-specialists in corpus linguistics, namely, second language learners and teachers. The MOOC course on Common Law, which was employed for the experiment presented here, is introduced by several YouTube tutorials, which explain briefly what the platform offers and how to exploit it to its fullest.

The MOOC course on common law used in this research could be deemed a corpus inasmuch as it contains a set of transcriptions of authentic legal English lectures given by Professor Adam Gearey, at the University of London for London Coursera. The prospective learners can watch the video of the lecture, which is also duly transcribed for them to read and work, having recourse to all the different functionalities which the platform offers. As regards the content of the lectures, they deal with various issues such as the history of common law, the structure of courts and tribunals in Great Britain, the sources of law, the principle of binding precedent or European law. The transcriptions themselves vary considerably in terms of their textual features, some of them belonging to the oral mode due to the presence of questions, inserts, pauses or simple syntactic structures, which denote the speaker’s intention to catch the listeners’ attention and to keep them engaged in the talk. In other cases, the lectures are often read by the speaker, being more formal as regards lexical choice, more syntactically complex and better planned and organised, as is typical of the written mode.

## A data-driven learning experiment in the legal english...

Amongst other functionalities, the FLAX system facilitates the retrieval of typical word and phrase usage samples by grouping linguistic data and sorting search results to show the most common patterns. It is capable of producing term lists, like the one illustrated by figure 1, which allow the user to search for the most relevant concordances associated with each of these.

academic Words	sort by frequency
legal	153
fundamental	75
link	52
legislation	49
role	38
individual	33
authority	29
define	27
couple	25
paragraph	24
similar	21
primarily	19
quote	15
somewhat	15
conventional	14
final	12
quotation	12
image	11
commission	11
restriction	10
civil	115
interpretation	73
theme	52
interpret	48
involve	38
source	32
structure	29
period	27
method	25
debate	23
function	21
evidence	18
policy	15
legislative	15
access	13
ensure	12
so-called	12
element	11
enable	10
affect	9
convention	115
obviously	73
process	52
institution	48
create	38
stress	32
tradition	29
contemporary	27
concept	25
integrity	23
justify	21
circumstance	18
chapter	15
grant	15
achieve	13
resolve	12
definition	12
tension	11
complex	10
enforce	9
principle	95
area	71
community	51
constitutional	46
context	35
lecture	32
constitution	29
major	27
impact	25
inconsistent	23
focus	20
prior	18
instance	15
labour	14
accurate	13
creation	12
precedence	12
prohibit	11
feature	10
precise	9
issue	84
previous	57
precedent	51
hierarchy	41
economic	34
domestic	32
require	29
distinction	26
section	24
consistent	23
approach	19
presumption	17
relevant	15
ultimately	14
coherent	13
ambiguity	12
commentator	12
specific	11
category	10
abstract	9

Figure 1. Legal term list (wordlist function on the menu)

Secondly, it automatically retrieves collocations and lexical bundles according to part-of-speech tags—for instance, all the adjectives associated with a particular noun—, as shown in figures 2 and 3. Learners can explore these elements by searching and browsing, and inspect them along with contextual information. The platform also presents them with general and academic English words, hyperlinked to their usage and collocates in authentic contexts.

**English Common Law MOOC (University of London with Coursera)**

About Collection Search Lectures Quizzes Extras Activities Collocations Wordlist LexicalBundles My Cherry Basket

**Browse Collocations in Collection**

a b c d e f g h i j k l m n o p q r s t u v w y Top 100

**9 collocation(s) associated with the word appellate**

Adjective (6) Verb (2) Noun + of (1)

- + appellate jurisdiction (8)
- + appellate courts (5)
- + appellate capacity (2)
- + appellate level (2)
- + appellate structure (1)
- + appellate committee (1)

Figure 2. Collocations of the term appellate.

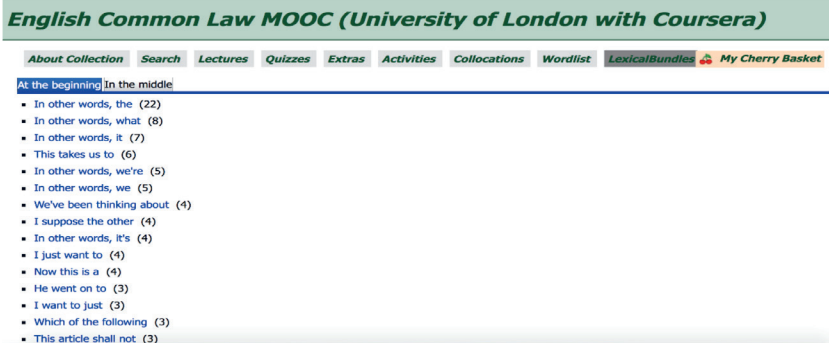


Figure 3. Lexical bundles.

44

One of the most useful functionalities offered by the FLAX is the possibility of exploiting term usage by working on the activities proposed in each section such as “completing collocations”, “word guessing” or “scrambled sentences”, amongst others, and also of consulting other contexts like Wikipedia by activating the “wikify” option, where the most salient terms are linked to their definition and related topics. The definition of the term *appellate jurisdiction*, which appears in a green text frame, is retrieved by the system from Wikipedia. Those elements which are “wikified” are highlighted in blue, allowing the user to see the definition and related topics comprised in it. However, due to the instructions given to the experimental group, as shown in the methodology section, and with the aim of not letting other sources “contaminate” the process, students were instructed not to activate this option during the present experiment.

Thus, all the different functionalities offered by the FLAX platform might make it a suitable tool to be employed in corpus-based language instruction as its design addresses some of the challenges Ädel (2010) poses within the field. Firstly, Ädel complains about the lack of available academic corpora (which is particularly remarkable in the legal field) and the growing demand for this kind of materials, which the FLAX offers online and exploits through all the possibilities described above. As regards hands-on work with corpora, Ädel also detects some problem areas that corpus-based instruction must cope with, some of which, in our view, could be tackled if working with a system like the FLAX.

Like most of its critics (Widdowson 2000; Flowerdew 2009), Ädel refers to the decontextualised samples obtained when exploring corpora in the language classroom and to the corpus, as looking like a maze when presented to students, who often get



lost in the vast amount of information retrieved by corpus tools and can even get drowned in data, to use Ädel's words. There is a need to control such a large amount of input, a challenge that the FLAX addresses (at least on a lexical level) firstly, by filtering results and offering different options such as highlighting only terms, academic or general vocabulary and expanding their contextual information. Such expansion is carried out by connecting the selected vocabulary to the web through the wikify option. This facilitates, on the one hand, understanding and, on the other hand, presents the terms in various contexts acting as reference for later use.

Concerning the challenges of interpretation and evaluation of the information retrieved from corpora, as presented by Ädel (2010), it could also be argued that the FLAX partially addresses such challenges insofar as it selects the most relevant terms, collocational patterns or lexical bundles in a text collection and allows the learner to explore their contexts of usage. The system also guides the students through different activities to exploit term and collocation usage and, in a way, contribute to their acquisition.

### **3.2. Description of the Experiment**

The experiment presented herein could not be regarded as a DDL experiment proper but rather as a corpus-based self-learning experience which attempts to test the effectiveness of an online learning platform, the FLAX, used as a support in the legal English classroom. One of its major aims is to try to quantitatively determine the usefulness and effectiveness of employing the FLAX in the teaching of legal English. To that end, a group of 105 students in the fourth year of the Translation Degree at the University of Murcia (Spain) studying a legal English course were selected as informants. All the students' linguistic competence level complied with the CEFR requirements for the B2 level, having passed general English exams B1 and B2 prior to studying legal English.

Our initial intention was to incorporate the FLAX as part of the course methodology itself, trying not to alter the original syllabus of the subject in its essence. In order to do so, the informants, who had to write an essay on a given set of legal English topics —defined by the subject instructor— as part of their final assessment, were divided into two groups. The experimental group (34 informants organised in 8 subgroups) were requested to only consult the FLAX website as the single source of information to draft their essays. The remaining 71 students (divided into 16 different groups) would act as the control group, following the usual working method for the design and drafting of their work, that is, they could employ any information source available without any limitation or previous instruction. The groups were not balanced because the FLAX course lessons included in the MOOC course on English Common Law did not cover some of the topics comprised in

the syllabus of the subject (designed before this experiment took place), and only 8 out of 24 topics coincided with the ones listed on the online learning platform.<sup>3</sup>

The informants employing the FLAX (the experimental group) were instructed on its use in one session of one hour, where they were requested to follow the video tutorials provided in the legal English section of the website.<sup>4</sup> This would imply, not only watching the videos on various legal English topics and reading their transcriptions, but also using the functionalities present in every lesson as well as the language activities described above. The informants were further instructed to abstain from consulting any reference outside the platform, being constrained to use the sources supplied by it, whereas the control group was given the liberty to resort to any kind of source and/or reference such as related bibliography or internet websites dealing with the subjects involved. After following all the steps described in the tutorials, all the members of the experimental group would start writing their essays trying to incorporate all the relevant information and the specialised terminology required in each case.

### 3.3. Learner Corpora Description

46

Once the essays were finished, they were gathered forming two small learner corpora whose size differed considerably for the reasons explained above (see section 3.1). The FLAX-based corpus contained 34,647 tokens,<sup>5</sup> while those texts not based on it amounted to 108,681. The extension of the texts in each corpus ranged from 2,356 tokens to 10,908. On the whole, those texts which were not based on the FLAX tended to be longer, including 6,393 tokens on average, as opposed to those based on the FLAX, containing 4,330. The fact that there were many more data available for those informants using the internet or other information sources might account for this noticeable difference.

Both corpora were processed automatically using Scott's (2008a) *Wordsmith Tools* with the purpose of extracting information tending to hint at the suitability of the FLAX as an experimental learning method, as opposed to the usual working method used by the subjects in the control group. In the first place, the texts were analysed applying Corpus Linguistics techniques for the exploration of the lexical level of the language, namely, lexical diversity, specialised term usage and distribution and lexical fundamentality.

Additionally, and with the aim of revealing the interpersonality traits in the texts of either corpus, a thorough analysis of the MD markers present in them was undertaken. Such an analysis was deemed necessary to go beyond the lexical choices made by the informants to ascertain whether the stance taken by the authors of the essays in either group towards the propositional content of their work bore relevant differences attributable to the use of FLAX or non-FLAX materials.

## 4. Results and discussion

### 4.1. Lexical Analysis

#### 4.1.1. *Lexical Diversity*

One of the possibilities offered by Scott's *Wordsmith 5.0* is to compute the type/token ratio in a corpus, that is, the proportion existing between a word (type) and the number of occurrences of that same word in the corpus (token). When the size of the texts in each corpus is different, Scott (2008b: 221) recommends applying the standardised type/token ratio (STTR), which is calculated for the first 1,000 tokens in each text, since, when text length varies, the results may also differ considerably, as is the case in the present study. A corpus with a high STTR would contain a higher number of types per token than one with a lower ratio, consequently, the breadth of its vocabulary would necessarily be greater.

Regarding the corpora under examination, the one not based on the FLAX displayed a higher STTR, reaching 37.63 as opposed to the FLAX-based text collections, over 2 points below (35.3). As already stated, a higher STTR would necessarily imply greater lexical diversity; therefore, although the difference is not substantial, those texts written using the internet and other bibliographic sources displayed greater vocabulary breadth, whereas the lexicon in the FLAX-based corpus tends to be more repetitive, according to the figures.

This fact is directly related to the observations presented in section 4.1.4. on lexical fundamentality, which refers to the amount of general vocabulary found in both text collections. The results presented in 4.1.4. reinforce our perception of the smaller vocabulary breadth of the texts in the FLAX corpus as measured by STTR. These texts also present higher frequency of general vocabulary than those containing information from various sources other than the FLAX. As a consequence, it could be stated that, in spite of the greater proportion of specialised terms and their better distribution in the FLAX text collection (as illustrated below), these texts also display a poorer vocabulary when it comes to expressing non-specialised ideas and concepts.

#### 4.1.2. *Specialised Term Usage*

On a lexical level, one of the parameters that was measured was term usage. The relevance of terms in academic texts is fundamental as they could be regarded as conceptual vehicles which can be employed to transmit specialised knowledge amongst scientists, researchers, professionals or language learners, as is the case. In Kit and Liu's words, terms are "linguistic representations of domain-specific key concepts in a subject field that crystallise our expert knowledge in that subject" (Kit and Liu 2008: 204).

In order to quantify term usage, both corpora were analysed using Scott's (2008a) *Keywords* functionality included in the *Wordsmith 5.0* software package, a powerful corpus analysis tool which, according to Marín (2014a), turned out to be one of the most efficient in the extraction of legal terms from an 8.85 million-word legal corpus, the *BLaRC* (the *British Law Report Corpus*), reaching a peak of precision of 85% for the top 200 candidate terms identified by it.<sup>6</sup> Following Scott (2008b: 104), a word is key “if it is unusually frequent (or unusually infrequent) in comparison with what one would expect on the basis of the larger word-lists”, that is to say, its degree of specificity could be related to its keyness given its statistical behaviour both in the general and the specialised fields.

*Keywords* managed to mine 349.2 specialised terms from the learner corpus based on the FLAX and 309.1 from those texts not using the FLAX after normalisation. The difference in size between both corpora led to the normalisation of the data obtained, which consisted in dividing the total number of terms extracted from each corpus by the number of tokens in them. Subsequently, the figures were multiplied by 100,000 to avoid an excessive amount of decimals. In order for the list of candidate terms produced by *Keywords* to be validated, it was compared against a glossary of 10,088 legal terms<sup>7</sup> so that every time a candidate term was found in the glossary, it was confirmed as such.

48

In spite of the similar number of terms in each corpus, their proportion with respect to the whole type list was three times as high (10.32%) for the FLAX-based corpus as for the non-FLAX-based one (3.82%). It could, therefore, be argued that those students using the FLAX as an information source for the drafting of their essays, showed greater command in the use of legal terms than those who did not.

The observed data related to the proportion and average frequency of specialised terms in both corpora were also scrutinised from the perspective of inference statistics.<sup>8</sup> Inference statistics, amongst other possibilities, allows linguists to make generalisations on the language based on the observations of a given sample. It “pertains to the need to generalise from a finite sample of language data to a theoretical infinite amount of text” (Baroni and Evert 2009: 779). Using the average frequency of specialised terms in both corpora and the number of tokens in each of them, the probability for these to occur in a hypothetical total population of similar texts also indicated that it was higher for the FLAX collection obtaining a frequency estimate of 11.77% as against 4.81% for the non-FLAX set.

Table 1 displays the top 20 legal terms extracted from both sets of texts using Scott's keywords.

NON FLAX-BASED TEXTS (control group)		FLAX-BASED TEXTS (experimental group)	
TERMS	KEY-NESS	TERMS	KEYNESS
Law	7833.93	Law	3584.82
Contract	3050.50	Rights	1586.21
Legal	2839.63	Court	1378.70
Civil	2493.52	Precedent	1187.08
Attorney	1904.97	Case	702.25
Court	1577.73	Sovereignty	641.44
Criminal	1361.42	Statutes	468.01
Offence	1316.28	Act	467.83
Party	1266.36	Decisions	429.75
Custody	796.83	Convention	372.03
Testator	649.71	Appeal	337.26
Property	600.88	Legislation	227.05
Probate	581.39	Rule	219.71
Contractual	531.89	Civil	210.75
Power	523.25	Constitution	210.63
Legislation	509.75	Power	201.92
Arbitration	485.55	Interpretation	197.23
Act	432.32	Binding	184.37
Notary	426.29	Judicial	179.32
Agreement	422.80	Jurisdiction	158.02

Table 1. Top 20 legal terms

One of the major conclusions to be drawn with respect to the top 20 terms identified in both corpora pertains to the nature of such terms. For instance, a term like *attorney* could not be found in the FLAX-based text collection because that collection only includes British texts and *attorney* is a legal term from the American system. Furthermore, most of the terms in the FLAX-based list refer to the sources of law and the norm itself (*precedent, act, statute, constitution, legislation, convention*) and also to their procedural application (*court, interpretation, judicial, appeal, binding*), whereas the non-FLAX term sample displays greater heterogeneity since, although it contains some of these terms (*act, legislation*), it does not refer to the major source of law *par excellence*: case law. In

fact, it gathers terms from various legal areas, mainly contract law (*contract, contractual, agreement, party, arbitration*), and also property and its management (*probate, property, testator, notary*).

In spite of such difference, the specificity level of the terms identified in both corpora differs considerably. Using *Keywords* as the tool to mine the most relevant terms by comparison with a general English corpus, LACELL, of 21 million tokens, it was found that the terms in the non-FLAX corpus displayed an average keyness value of 179.33, whereas those using the FLAX as a resource stand 20 points below, at 156.49 keyness.

However, specialised terms represented 10.32% of the whole type list in the FLAX corpus as opposed to the non-FLAX text collection, where the percentage of terms identified is three times lower, that is, 3.82%, as shown above. It could be assumed that the experimental groups used the terminology more consistently than the control group although the latter employed more specific terms than those used by the former. As examples of usage by both groups:

- (1) In a will, the testator or testatrix appoints another person (called the executor) as responsible of the administration and distribution of his/her possessions among his/her inheritors or beneficiaries (Non-FLAX).
- (2) A.D.R consists of choosing a judge called arbitrator that, after examining the different positions of the parties, issues a binding decision called arbitration (Non-FLAX).
- (3) The term binding precedent is the opposite idea to persuasive precedent, which is not binding (FLAX).
- (4) The parliament (...) creates supreme law (statutes), which will override inconsistent case law and reflect the sovereignty and legitimacy of parliament (FLAX).

#### 4.1.3. Term Distribution

The distribution of terms within a learner specialised corpus is also a relevant piece of data which can reveal information on the learners' knowledge of the terminology and their capacity to employ it in a wider set of contexts. As a matter of fact, the word *distribution* is used in this study to refer to the amount of texts in a corpus where a term can be found: it is expressed in percentages to respect proportionality. Therefore, the better distributed a term is, the more relevant it might be to the corpus.

Distribution, or text range, can be computed automatically using the *Wordlist* software included in the *Wordsmith* package (Scott 2008a), as well as the type/token ratio, which are provided within the general statistics. In this particular

case, it must be taken into account that the texts in the corpus deal with different legal areas (family law, European law, civil law or common law, amongst other), thus, except for the most general terms covering a wider range of topics, the majority of the terms extracted from each corpus should rather be restricted to their specific areas. Nevertheless, the figures show that the specialised terms used by the experimental group are much better distributed, occurring in 48.49% of the corpus texts on average, whereas the control group average value for this parameter is noticeably lower, 29.54%. These percentages reflect mean values, that is to say, terms like *convention* or *ruling* appear in 100% of the texts in the FLAX corpus while *override* or *injunction* are only employed in 34% and 25% of the texts respectively. If we consider the whole list of specialised terms obtained from both text collections, they are better distributed in the experimental group, whose term list covers almost half of the texts included in the corpus.

Even so, these figures can be read in different ways. On the one hand, it appears that the informants in the experimental group may have a wider knowledge of the terminology, as they are capable of using terms which are not only related to the legal area they have researched but also to other areas present in other corpus texts, given the high average distribution percentage obtained (48.49%). In fact, the legal terms identified in the FLAX corpus can be found in almost half of its texts. Terms such as *law*, *court*, *case*, *rule* or *convention*, appear in the whole of the text collection, whereas others like *injunction*, *litigant* or *jurist* are limited to just one of the texts, due to their more specific character.

On the other hand, the lower distribution percentages computed for the control group, at almost 20 points below the experimental one, might well be related to the learners' more limited knowledge of the terminology, although it may also be associated with the more specialised meaning of the terms used in this text collection, being found in fewer texts in the corpus. This hypothesis might be supported by the average keyness value of the terms in both lists, 20 points higher for the control group, which could be indicative of the greater specificity of the terms found in the non-FLAX texts. Either way, in order to confirm this perception, based on the data obtained automatically from both text collections, a manual scrutiny of the texts included in each corpus would be required to complement this quantitative analysis.

#### 4.1.4. *Lexical Fundamentality*<sup>9</sup>

Processing both corpora with the software *Range* (Heatley and Nation 1996) could also provide an insight into the lexicon of both text collections. The version employed in this study is the one offering the possibility of processing

the texts in a corpus in comparison with the most frequent 3,000 words found in the *British National Corpus*<sup>10</sup> (BNC), a general English corpus of 100 million words. This software allows the user to calculate text range, that is, the percentage of running words in a corpus covered by those 3,000 words which are arranged in sets of 1,000 according to their frequency in it. The figures below were obtained by comparison with the first list only of the most frequent 1,000 words in this general corpus. Words such as *and*, *baby*, *because*, *hate*, *the* or *then* could be found within that list. As a consequence, the higher the text range percentage obtained after processing a corpus, the more fundamental the lexicon in that corpus. On the contrary, if the percentage of tokens covered by these lists was lower, the vocabulary in a corpus would necessarily be more specialised, or at least less basic.

Concerning our two corpora, lexical fundamentality was computed automatically by processing them with *Range*. The highest percentages were assigned to the FLAX corpus, reaching 79.39% text range, while 20.61% of the tokens in that corpus could not be found in the BNC lists of the most frequent 1,000 types in it. In contrast, only 66.73% of the types in the non-FLAX corpus overlapped with the ones on the BNC lists. These percentages indicate that the former corpus displays greater lexical fundamentality than the latter, that is, it contains a higher number of tokens present in the lists of the most frequent/basic types of English used as reference for their processing.

This finding might contradict the results discussed in section 4.1.2., where it was observed that the ratio of terms per token was higher in the FLAX corpus in spite of its lexicon seeming more fundamental or basic, as illustrated by the percentages above. Nevertheless, it could also be argued that legal terms such as *case*, *rule* or *rights*, in spite of being considered as specialised terms, could be found amongst the most frequent 3,000 types of English. Their sub-technical character accounts for this fact, since they are shared both by the legal and the general contexts. On the contrary, the use of terms like *testator*, *probate* or *arbitration*, included in the top 20 legal terms extracted from the non-FLAX texts, could also explain this fact. They are much more specific and tend to be employed in fewer texts, hence the lower distribution values discussed in the previous section.

Even so, the lower Standardised Type Token Ratio (STTR) associated with the texts in the FLAX collection, may also reinforce our perception that, although more specialised in the way they refer to legal concepts (judging by term ratio and distribution figures), the texts in the FLAX corpus display a tendency on the part of the authors to use more general vocabulary which also, in general, tends to be slightly less varied.



#### 4.2. Analysis of Metadiscourse Markers: Results and Discussion

As was anticipated at the end of Section 3, a study of the presence of MD (metadiscourse) markers of the textual and interactional kind was also implemented using Scott's WordSmith 5.0 tool with the aim of studying their statistical behaviour in both the FLAX and non-FLAX corpora. The goal was to reveal differences in the way in which propositional content was presented as regards writers' engagement and stance, as specific samples of the RA (research article) genre.<sup>11</sup>

As shown in Figure 4, the results indicate that, in both corpora, the overall number of textual MD markers was much higher than the set of interactional ones. They also attest that these textual markers were more frequently employed in the control group than in the experimental one (554.61 against 452.83, respectively<sup>12</sup>). Contrarily, interactional markers occurred more frequently in the FLAX group, displaying 154.06 frequency as opposed to 143.24 (non-FLAX texts).

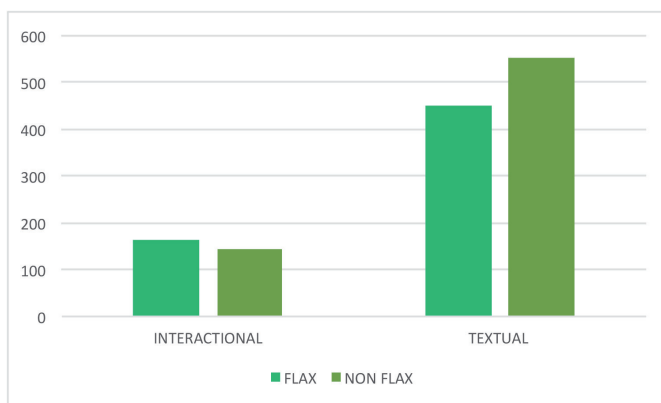


Figure 4. Metadiscourse markers in non-FLAX and FLAX-based corpora

The main reason for the greater number of textual MD markers might indicate an overall preference on the part of the informants to convey propositional content in an orderly manner, rather than engaging with the prospective readers through evaluation, appraisal and other affective resources. After all, the informants are a group of highly proficient undergraduate translation students who might lack enough self-confidence in the area of academic writing to mobilise the empathy of the prospective readers, focusing upon achieving grammar correctness and adequacy instead.

More specifically, as reflected in Table 2, transition/logical markers are the most numerous ones in either corpus (with 403 occurrences in the non-FLAX set and 384 in the FLAX-based set). This finding is in line with the claim made by Hempel and Degand (2008) concerning the importance of textual markers used in various texts, these resources being the authors' conscious stratagem in constructing the propositional content which they aim to convey to the addressee. In this sense, 'and' is, by far, the most recurrent connector in either corpora, followed by 'or'. This data might indicate that there is an overall marked preference for linking ideas through additive markers, and, for second choice, the use of adversative markers to construct arguments (Dafouz 2008).

TEXTUAL MARKERS		
TYPES	NON-FLAX CORPUS (norm. freq.)	FLAX CORPUS (norm. freq.)
<b>TRANSITION/LOGICAL MARKERS</b>		
and	249.17	246.77
furthermore	1.29	1.15
additionally	0.55	0
or	104.43	53.11
but	16.47	27.71
however	7.45	10.68
nevertheless	2.85	6.64
so	8.83	19.92
therefore	3.40	7.22
finally	2.85	2.60
moreover	1.20	1.73
hence	0.28	1.44
thus	2.48	4.04
in addition	1.29	0.87
in summary	0.00	0.29
in conclusion	0.09	0.29
what is more	0.09	0.00
concluding	0.37	0
<b>SUBTOTAL</b>	<b>403.09</b>	<b>384.46</b>

A data-driven learning experiment in the legal english...

ENDOPHORIC MARKERS		
noted/see above/below	1.74	0.86
see fig	0.09	0
in section X	0.18	2.3
<b>SUBTOTAL</b>	<b>2.01</b>	<b>3.16</b>
FRAME MARKERS		
in the first place	0.09	0.29
firstly	0.64	1.15
as stated in	0.28	0.00
as for	0.74	0.00
as regards	0.00	0.29
thirdly	0.18	0.29
secondly	0.74	0.87
regarding	4.05	2.31
concerning	1.38	1.15
<b>SUBTOTAL</b>	<b>8.1</b>	<b>6.35</b>
CODE GLOSSES		
that is	2.58	3.17
in other words	0.18	0.00
explicitly	0.18	0
specifically	0.83	0.29
—	0.92	0.00
()	119.52	25.40
colon	8.37	18.76
namely	0.28	0.29
<b>SUBTOTAL</b>	<b>132.86</b>	<b>47.91</b>
EVIDENTIALS		
according to X	8	6.63
X states/says	0.55	4.32
<b>SUBTOTAL</b>	<b>8.55</b>	<b>10.95</b>
<b>TOTAL</b>	<b>554.61</b>	<b>452.83</b>

Table 2. Textual markers in non-FLAX and FLAX-based corpora

The use of other, more sophisticated kinds of connectors is negligible by comparison in either corpus, in tune with Moreno's (2004: 21) findings on the dearth of textual indicators in Spanish academic corpora, or their comparative scarcity with respect to English academic writing (Mur Dueñas 2011: 3071). As examples of 'and' in each corpus:

- (5) It comprises the rule by which a court hears and determines what happens in civil lawsuits (Non FLAX).
- (6) Defamation: it occurs when the defendant communicates untruthful information about the plaintiff and it hurts the plaintiff's reputation (FLAX).

The next group with the most markers (132.86 and 47.91) is code glosses. Glosses are used by writers to ensure the readers understand the meanings of specific elements, phrases, or idioms. Again, this kind of explanatory device is markedly more present in the control group (132.86), almost exclusively in the shape of parentheses as a means to expand, define or delimit the propositional content. It would suggest that the informants in the non-FLAX group are aware of the complexity of the subject they are dealing with, providing their audience with a number of explicit reading prompts as well as more examples, in the attempt to render their explanations clearer. On the other hand, the amount of code glosses employed by the experimental group is much smaller (47.91), the occurrences taking place, as in the control group, mostly through parentheses (25.4), but also with a relatively high number of colons (18.76). As examples of group glosses other than parentheses and colons:

- (7) (...) which are not considered as crimes nor breaches of contract, that is, torts. (Non FLAX)
- (8) In other words, they tried to make a case that would not be a precedent (sic). (FLAX)

Frame markers are comparatively less present in either group, even if again they are more abundant in the control one, with 8.1 and 6.35 occurrences respectively. Frames organise sequences, label text stages, announce topic goals and indicate topic shifts. The scarcity in both corpora (6.35 for the FLAX-based texts and 8.1 for the non-FLAX ones) might mark the same dearth of sophistication in academic writing that was pointed out when discussing the simplicity of the logical connectors deployed by the two groups of informants. Finally, evidentials are used to inform readers about who has said or written a given idea or comment. Mainly, they are used by way of testimonials that give weight to the supposed value of propositional content reflected by the authors, sustaining and validating their theses. The presence of evidentials is also scarce in either corpus, even though these, together with endophoric markers, are more numerous –albeit marginally–

in the FLAX texts. Evidentials display 10.95 occurrences in the FLAX group, compared to 8.55 occurrences in the non-FLAX group, the writers in the former seemingly exhibiting greater awareness of the need to establish their credibility through the knowledge of the ‘right’ texts.

As far as intratextual references (or endophoric MD markers) are concerned, their appearance is also scarce, but slightly more frequent in the FLAX-based corpus, with 3.16 occurrences, as opposed to 2.01 in the control group. As Heng and Tan (2010) discovered, the use of endophorics –used to support the argument by convincing readers of the validity of the argument– could be closely linked to the use of citation as a persuasive strategy in the crafting of academic writing. This affirmation would be in line with our conclusions below, pointing to the fact that the FLAX corpus could show a subtly higher degree of sophistication and capacity of persuasion if compared with the resources used by the non-FLAX group.

The results of the scrutiny of interactional markers in the corpora are reflected in Table 3 below.

INTERACTIONAL MARKERS		
TYPES	NON-FLAX (norm. freq.)	FLAX (norm. freq.)
<b>HEDGES</b>		
May	16.47	16.16
Might	1.75	3.17
Must	17.21	11.55
Can	35.06	43.29
Could	3.50	16.45
Would	5.89	27.42
Probably	0.28	0.87
Perhaps	0.09	0.87
Maybe	0.09	0
<b>SUBTOTAL</b>	<b>80.34</b>	<b>119.78</b>
<b>BOOSTERS</b>		
Clearly	1.10	1.15
Certainly	0	1.15
<b>SUBTOTAL</b>	<b>1.10</b>	<b>2.3</b>

ATTITUDE MARKERS		
need to	3.12	3.46
we think	0	0.57
I think	0	0.28
have to	4.50	10.39
Unfortunately	0	0.86
<b>SUBTOTAL</b>	<b>7.62</b>	<b>15.56</b>
ENGAGEMENT MARKERS		
consider that	0	0.86
<b>SUBTOTAL</b>	<b>0</b>	<b>0.86</b>
REFERENCES TO SELF		
I	9.01	4.61
Me	2.02	0.57
us	4.60	5.48
our	2.76	4.90
mine	0.36	0
<b>SUBTOTAL</b>	<b>18.75</b>	<b>15.56</b>
<b>TOTAL</b>	<b>143.24</b>	<b>154.06</b>

Table 3. Interactional markers in non-FLAX and FLAX-based corpora

As we anticipated at the beginning of this section, these markers occur less often in the texts under study, probably on account of a reluctance on the part of the informants to appraise the propositional content of the text. This result also agrees with Mur Dueñas's (2011: 3075) findings in a corpus of Spanish research articles, where she shows that Spanish writers tend to establish a smaller degree of interaction with their addressees than English ones do. Also, as beginners in the drafting of academic texts, the informants might be reluctant to show complicity with the reader, favouring the use of textual markers that organise the discourse in a more conventional way from an academic perspective instead.

Nevertheless, interactional markers are rather more dominant in the FLAX than in the non-FLAX corpus (154.06 against 143.24, respectively), pointing to the possibility that the experimental group might be comparatively more willing to interact with their readership and engage with it. Still, the thesis hinted at above, that, in general, informants in both groups might be more 'academically conservative', would be reinforced by the high presence of hedges in both corpora (119.78 in the FLAX collection against 80.34 in the non-FLAX set). Hedges –

mainly introduced by auxiliary ‘can’ in both corpora– are a usual device deployed by academic writers, since they can “anticipate possible opposition to their claims (by expressing statements with precision but also with caution and modesty), while simultaneously, enabling the reader to follow the writer’s stance without the writer appearing too assertive” (Dafouz 2008: 107), for instance:

- (9) As long as possible, we should transfer it, although it can be translated in some cases as “fideicomiso” (Non-FLAX).
- (10) The term Common Law can first of all be understood as the law imposed on the institutions of the Anglo Saxon England (FLAX).

In academic texts, they are normally counterbalanced by boosters, but the appearance of these is residuary in both our corpora, which strengthens our previous assertions. Additionally, the absence of engagement markers would again confirm the lack of commitment on the part of the writers both in the FLAX and non-FLAX groups.

Attitude markers, in turn, are also scarce, if somewhat more present in the FLAX group, but mostly through the modal auxiliary ‘have to’. Finally, the figures obtained account for similar results in the area of self-mentions, which are the only MD markers which the control group uses more often than the experimental one, although marginally (18.75 in the non-FLAX texts against 15.56 in the FLAX corpus). This is achieved mainly through the use of the first person singular. In both corpora, the first person plural pronoun is used to inform the writers of their intention, such as ‘we will now deal with’, ‘we will then present’ and ‘we will include’, thus indicating authorial presence, not only of the individual informant, but also of the working team as a whole. Below are some usage examples:

- (11) Henceforth, we will focus on civil law from Common law and its division (Non FLAX)
- (12) However, we must not forget that history is fuel to the future and that our current idea of due process is (...) (FLAX).

Like specialised terms, inference statistics confirms our perception about the two major groups of MD markers. While textual markers would represent a comparable proportion of texts within a hypothetical population of such linguistic units, that is, 1.44% and 1.46% for the FLAX and non-FLAX corpora respectively, the greatest difference would be found amongst interactional markers, obtaining 15.4% frequency estimate for the former and 14.32% for the latter.

In sum, MD markers are present in the corpora under study, as specific samples of RAs, where students are initiated in the writing of academic genres. Nevertheless, they occur in the most conventional ways, i.e. through the use of textual indicators aimed at arranging, organising and ‘tidying up’ the propositional content in the

texts. Both corpora, mainly the non-FLAX one –with its abundance of logical connectors–, are conventionally constructed inasmuch as they fit the impersonality and detachment that traditionally surround academic texts. Nevertheless, if persuasion is also a desirable element in this kind of texts, it is not to be found in either of the corpora under analysis. Certain differences in engagement between the control and experimental groups are observed, but these are not significant, since both corpora adopt predictable devices, mainly logical transitions and, chiefly in the control group, explanatory glosses. In the area of stance, i.e. of interactional markers, it is the FLAX group that shows a greater degree of sophistication. Within this category, hedges –at a greater distance from other groups, with the slight exception of self-mentions–, are the most favoured MD markers, which, again, could indicate a relatively primitive state of affairs in the informants’ writing abilities.

## 5. Conclusion

60

This research has attempted to quantify the usefulness of corpus-based materials used as support to the legal English classroom. One of the key factors which motivated it was the fact that DDL experiments in this ESAP variety are scarce, leaving room for greater experimentation and speculation about the benefits of implementing such methods in legal English teaching which, to the best of our knowledge, remains underexplored in the literature.

To that end, the FLAX, an online language learning platform offering a course which contains a corpus of university lectures on legal issues, was used as part of an experiment where two groups of informants were instructed to write academic essays on legal topics. The FLAX was used by the experimental group as their only source of information while the control group could consult any reference at hand for the same task. As already stated, the FLAX addresses some of the challenges posed by Ädel (2010) which remain to be met by DDL methodologies. On the one hand, the scarcity of academic corpora available (which is particularly remarkable in the legal field) is a major concern to this author. In this respect, the FLAX offers free access to legal corpora, which are exploited through the proposal of language activities and other functionalities. In addition, Ädel (2010) also deems raw corpus data to be a “maze” which learners have to go through often getting “drowned” by the vast amount of data generated by concordancers. In this respect, the FLAX filters the information retrieved from corpora through term/vocabulary lists which are offered in context and linked to other information sources.



As regards the two research questions formulated in the introduction: firstly, we wondered whether the FLAX would positively influence the usage of specialised legal terminology. The answer to this question would be affirmative since, on a lexical level, after processing the two learner corpora gathered for this study, the figures indicate that the experimental group used the specialised terminology better than the control group, utilising 10.32% specialised terms for the expression of technical concepts as opposed to the non-FLAX corpus, where the presence of legal terminology was three times lower. Term distribution was also higher in the FLAX corpus, standing at 20 points above the same value for the control group (28%). Nonetheless, the lexicon employed by the experimental group appeared to be poorer, as attested by the standardised type/token ratio values yielded after processing both corpora. Although the difference was not substantial, the proportion of different types was greater in the non-FLAX corpus and hence the diversity of its lexicon. Likewise, it was noted that the lexicon of the FLAX corpus tended to be more basic than the corpus obtained from the experimental group, as 79.39% of the types found in it overlapped with the list of the 1,000 most frequently used words taken from the *British National Corpus*. Whether in fact this turns out to be a disadvantage of this teaching-learning method would require further research.

61

The second research question posed in the introduction could also be answered affirmatively. In the first place, as has been illustrated throughout section 4.2, corpus linguistics could throw light on the decisions made by second language learners on a pragmatic level in the deployment of metadiscourse markers. As a matter of fact, the use of these elements in both our corpora showed slight differences. This was shown by the way in which textual markers were employed by the informants, mainly logical transitions and glosses, which were more abundant in the text collection produced by the control group. On the other hand, interactive markers showed a lesser presence in both our corpora, probably due to reluctance on the part of the informants –as we may recall, English non-native undergraduate students– to appraise the propositional content of the text.

Nevertheless, the greater deployment of persuasion in the shape of interactive markers in the FLAX group indicated that the experimental group was comparatively more willing to interact with the readership and engage with it. It could be argued that such willingness might be a consequence of the text genres found in the FLAX platform. The online texts accompanying the videos of the lectures transcribe Professor Gearey's lessons literally, presenting certain features of oral language, itself necessarily of an interactive nature. However, many of the lectures are read by the speaker and also present clear features of academic writing, and so, it cannot be stated for certain that there is a direct relation between the texts written by the

experimental group and the textual genres the transcriptions might adhere to.

On the whole, it could not be categorically stated that the use of the FLAX benefited its users dramatically, although the analysis above illustrates a tendency on the part of the experimental group (only using the FLAX as their information source) towards utilising the terminology more consistently and employing MD markers more often, albeit marginally, for the expression of persuasion. Even so, further research would be needed along these lines to reach sounder conclusions and reinforce our initial perceptions.

## Notes

---

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

<sup>2</sup> Antconc (Anthony 2011) or the more sophisticated Wordsmith tools (Scott 2008), which would necessarily require training prior to actually engaging into the learning process itself.

<sup>3</sup> FLAX was not designed *ad hoc* to be tested in this translation course but rather incorporated as part of the experiment *a posteriori*.

<sup>4</sup> See: <http://flax.nzdl.org/greenstone3/flax> (Law collections/English Common Law MOOC)

<sup>5</sup> The term “type” refers to every different word in a corpus, whereas “token” stands for the number of repetitions of the same word within it.

<sup>6</sup> This means that 85% out of 200 terms automatically identified by *Keywords* were confirmed as true terms after comparing them with a legal English glossary.

<sup>7</sup> This glossary was compiled by merging together and filtering three online legal glossaries found at: <http://www.legislation.gov.hk/eng/glossary/homeglos.htm>

<http://www.judiciary.gov.uk/glos>

<http://www.nolo.com/dictionary>

[http://sixthformlaw.info/03\\_dictionary/index.htm](http://sixthformlaw.info/03_dictionary/index.htm)

<sup>8</sup> The online frequency estimate calculator found on <http://sigil.collocations.de/wizard.html> was used to that end.

<sup>9</sup> This term has been taken from Ishikawa (2015), who also studies the presence of general vocabulary in the speeches and writings of Asian learners of ESL and refers to the proportion of general vocabulary found in corpora as lexical fundamentality.

<sup>10</sup> <http://www.natcorp.ox.ac.uk>

<sup>11</sup> A further study –but out of our scope– in line with the English-Spanish contrastive analyses performed by Moreno (2004) and Mur Dueñas (2011) would be interesting, taking into account the characteristics of the oral online corpus that the students departed from.

<sup>12</sup> The figures indicate normalised frequency owing to the different size of both corpora. See section 4.1.1. for further details on normalisation procedures.

## Works cited

---

- ÄDEL, Annelie. 2010. "Using corpora to teach academic writing". In Campoy Cubillo, M.C., B. Bellés Fortuño and M.L. Gea-Valor (eds.) *Corpus-based Approaches in English Language Teaching*. London: Continuum: 39-55.
- ANTHONY, Lawrence. 2011. *AntConc (Version 3.2.2)* [computer software]. Japan: University of Waseda.
- ASTON, Guy. 1997. "Involving learners in developing learning methods: Exploiting text corpora in self-access". In Benson, P. and P. Voller (eds.) *Autonomy and Independence in Language Learning*. London: Longman: 204-214.
- BARONI, Marco and Stefan EVERT. 2009. "Statistical methods for corpus exploitation". In Lüdeling, Anke and Merja Kytö (eds.) *Corpus Linguistics: An International Handbook, 2*. Berlin: Mouton de Gruyter: 777-802.
- BIBER, Douglas and Edward FINNEGAN. 1989. "Styles of stance in English: lexical and grammatical marking of evidentiality and affect". *Text* 9 (1): 93-124.
- BOULTON, Alex. 2010a. "Data-driven learning: Taking the computer out of the equation". *Language Learning* 60 (3): 534-572.
- BOULTON, Alex. 2010b. "Learning outcomes from corpus consultation". In Moreno Jaén, M., F. Serrano Valverde and M. Calzada Pérez (eds.) *Exploring New Paths in Language Pedagogy: Lexis and Corpus-Based Language Teaching*. London: Equinox: 129-144.
- BOULTON, Alex. 2011. "Data Driven Learning: the Perpetual Enigma". In Roszkowski, S. and B. Lewandowska-Tomaszczyk (eds.) *Explorations across Languages and Corpora*. Frankfurt: Peter Lang: 563-580.
- BOULTON, Alex. 2012. "Corpus consultation for ESP. A review of empirical research". In Boulton, A., S. Carter-Thomas and E. Rowley-Jolive (eds.) *Corpus-Informed Research and Learning in ESP. Issues and Applications*. Amsterdam: John Benjamins: 261-292.
- CLEREHAN, Rosemary, Giselle KETT and Renee GEDGE. 2003. "Web-based tools and instruction for developing students' written communication skills". *Proceedings of Exploring Educational Technologies*. 16-17 July, Melbourne: Monash University. <[http://www.monash.edu.au/groups/flt/eet/full\\_papers/clerehan.pdf](http://www.monash.edu.au/groups/flt/eet/full_papers/clerehan.pdf)> Accessed May 31, 2017.
- CURADO FUENTES, Alejandro. 2004. "The use of corpora and IT in evaluating oral task competence for tourism English". *CALICO Journal* 22 (1): 5-22.
- DAFOUZ, Emma. 2008. "The pragmatic role of textual and interpersonal metadiscourse markers in the construction and attainment of persuasion: A cross-linguistic study of newspaper discourse". *Journal of Pragmatics* 40: 95-113.
- FAN, May and Xu XUNFENG. 2002. "An evaluation of an online bilingual corpus for the self-learning of legal English". *System* 30 (1): 47-63.
- FLOWERDEW, Lynne. 2009. "Applying corpus linguistics to pedagogy: A critical evaluation". *International Journal of Corpus Linguistics* 14 (3): 393-417.
- HADLEY, Gregory. 2002. "Sensing the winds of change: An introduction to data-driven learning". *RELC Journal* 33 (2): 99-124.
- HAFNER, Christoph and Christopher CANDLIN. 2007. "Corpus tools as an affordance to learning in professional legal education". *Journal of English for Academic Purposes* 6 (4): 303-318.
- HEATLEY, Alex and Paul NATION. 1996. *Range. Computer software*. Wellington, New Zealand: Victoria University of Wellington. <<http://www.victoria.ac.nz/lals/about/staff/paul-nation>> Accessed July 4, 2017.
- HEMPEL, Susanne and Liesbeth DEGAND. 2008. "Sequencers in different text genres: Academic writing, journalese and fiction". *Journal of Pragmatics*, 40: 676-693.
- HENG, Chang Swee and Helen TAN. 2010. "Extracting and comparing the intricacies of metadiscourse in two written persuasive corpora". *International Journal of Education and Development using Information and Communication Technology (IJEDICT)* 6 (3): 124-146.

- HUNSTON, Susan. 2007. *Corpora in Applied Linguistics*. Cambridge: Cambridge U.P.
- HYLAND, Ken. 2005. *Metadiscourse: Exploring Interaction in Writing*. London: Continuum.
- HYLAND, Ken and Polly TSE. 2004. "Metadiscourse in academic writing: A reappraisal". *Applied Linguistics* 25 (2): 156-177.
- ISHIKAWA, Shin'ichiro. 2015. "Lexical Development in L2 English Learners' Speeches and Writings". *Procedia, Social and Behavioural Sciences* 198: 202-210.
- JOHNS, Tim. 1986. "Microconcord: A language-learner's research tool". *System* 14 (2): 151-162.
- JOHNS, Tim. 1991. "Should you be persuaded: Two examples of data-driven learning". *English Language Research Journal* 4: 1-16.
- JOHNS, Tim. 1997. "Contexts: The background, development and trialling of a concordance-based CALL program". In Wichmann, A., S. Fiegelston, T. McEnery and G. Knowles (eds.). *Teaching and Language Corpora*. London: Longman: 100-115.
- KIT, Chunyu and Xiaoyue LIU. 2008. "Measuring mono-word termhood by rank difference via corpus comparison". *Terminology* 14 (2): 204-229.
- MAO, Luming. 1993. "I conclude not: toward a pragmatic account of metadiscourse". *Rhetoric Review* 11 (2): 265-289.
- MARÍN, María José. 2014a. "Evaluation of five single-word term recognition methods on a legal corpus". *Corpora* 9 (1): 83-107.
- MARÍN, María José. 2014b. "A Proposal to Exploit Legal Term Repertoires Extracted Automatically from a Legal English Corpus". *Miscelánea: A Journal of English and American Studies* 49: 53-72.
- MARÍN, María José and Piedad FERNÁNDEZ TOLEDO. 2015. "The Influence of Cognates on the Acquisition of Legal Terminology: Help or Hindrance? A Corpus-based Study". *Procedia-Social and Behavioral Sciences* 198: 320-329.
- MARTIN, James R. and Peter WHITE. 2005. *The Language of Evaluation*. Hampshire: Palgrave MacMillan.
- McENERY, Tony and Andrew WILSON. 1996. *Corpus Linguistics*. Edinburgh: Edinburgh U.P.
- McENERY, Tony and Richard XIAO. 2010. "What corpora can offer in language teaching and learning". In Hinkel, E. (ed.). *Handbook of Research in Second Language Teaching and Learning*. London: Routledge: 364-380.
- MORENO, Ana I. 2004. "Retrospective labelling in premise-conclusion metatext: an English-Spanish contrastive study of research articles on business and economics". *Journal of English for Academic Purposes* 3: 321-339.
- MUR-DUEÑAS, Pilar. 2011. "An intercultural analysis of metadiscourse features in research articles written in English and in Spanish". *Journal of Pragmatics* 43: 3068-3079.
- SANCHO-GUINDA, Carmen and Ken HYLAND. 2012. "Introduction: a context-sensitive approach to stance and voice". In Hyland, Ken and Carmen Sancho-Guinda (eds.). *Stance and Voice in Written Academic Genres*. Basingstoke, UK: Palgrave Macmillan: 1-11.
- SCOTT, Mike. 2008a. *WordSmith Tools version 5* [computer software]. Liverpool: Lexical Analysis Software.
- SCOTT, Mike. 2008b. *WordSmith Tools Help*. Stroud: Lexical Analysis Software.
- SINCLAIR, John. 1991. *Corpus, Concordance, Collocation*. Oxford: Oxford U.P.
- SINCLAIR, John. 2003. *Reading Concordances: An Introduction*. London: Longman.
- TODD, Richard Watson. 2001. "Induction from self-selected concordances and self-correction". *System* 29 (1): 91-102.
- WIDDOWSON, Henry G. 2000. "The limitations of linguistics applied". *Applied Linguistics* 21 (1): 3-25.

Received: 24 November 2016  
Accepted: 4 July 2017