

# Data-driven learning of collocations by Chinese learners of English: A longitudinal perspective

Tanjun Liu<sup>1</sup> and Dana Gablasova<sup>2</sup>

<sup>1</sup> Department of Applied Linguistics, Xi'an Jiaotong-Liverpool University, Suzhou, China  
Department of Applied Linguistics, Xi'an Jiaotong-Liverpool University, Suzhou, China

<sup>2</sup> Department of Linguistics and English Language, Lancaster University, Lancaster, UK

## Abstract

Collocations, a crucial component of language competence, remain a challenge for L2 learners across all proficiency levels. While the data-driven learning (DDL) approach has shown great potential for collocation learning from a shorter-term perspective, this study investigates its effectiveness in the long term, examining both linguistic gains and changes in learners' confidence about which words collocate. The effect of DDL was compared with learning collocations using a corpus-based collocations dictionary and non-corpus-based tools. Learners' experience with the tools was also explored through a questionnaire. The study employed a quasi-experimental research design with 100 Chinese learners of English as participants in two experimental groups and a control group. A novel corpus tool, #LancsBox (Brezina et al., 2015), was used in the DDL approach to identify and visualise collocations. The results showed that the learners in the DDL group increased their collocation knowledge at the end of the treatment and retained the gains three months later. The learners also reported a significant increase in their confidence about which words collocate. Both changes were found to be more substantial than the effects of using the corpus-based collocations dictionary or other tools. As for their experience, learners reported satisfaction with using corpora in their writing and, importantly, continued with corpus consultation three months after the end of the intervention. The findings have implications for integrating corpus consultation into learning practice both inside and outside of the classroom, showing that with sufficient training, DDL can provide an effective method to learning complex linguistic features such as collocations.

**Keywords:** Data-driven learning, corpus-based learning, corpora, collocation, #LancsBox

## 1. Introduction

The knowledge of formulaic language lies at the core of fluent and appropriate language use, both in first (L1) and second language (L2) (Siyanova-Chanturia & Pellicer- Sánchez, 2019). Collocations, words combinations that systematically co-occur with each other, such as *heavy rain*, *take advantage* and *expand rapidly*, are a key component of formulaic language and vocabulary knowledge, with a direct impact on learners' ability to communicate effectively in their L2 and to be perceived as competent users of the target language in different domains such as academic settings (Durrant, 2019). However, evidence from research and pedagogy indicates that learning collocations in English tends to be difficult for L2 learners at all levels of proficiency, including advanced L2 users (Boers et al., 2014; Chen, 2019; Gablasova et al., 2017; Siyanova-Chanturia, 2015). This difficulty has been attributed both to the linguistic complexity of collocations as well as to the pedagogical approaches adopted when teaching them. Research evaluating the effect of different pedagogical practices has shown that many traditional methods of teaching vocabulary do not seem to result in significant gains in terms of collocation knowledge (Pellicer-Sánchez & Boers, 2019).

This study seeks to contribute to our understanding of effective ways of teaching collocations in English by evaluating a corpus-based approach, referred to as data-driven learning (DDL), to teaching these formulaic units. DDL uses corpora, large electronic databases of language samples, to expose learners to authentic language use and to draw their attention to the target linguistic features (Boulton & Cobb, 2017). Since its inception, DDL has provided a valuable addition to other methods of language teaching and proved to be particularly effective at teaching complex linguistic concepts (cf. Chen & Flowerdew, 2018; Lee et al., 2019), where other teaching approaches have only a limited success (e.g., Vyatkina, 2016a). In particular, this study investigates the effectiveness of DDL for teaching collocations to advanced Chinese learners of English in an academic setting with the aid of a new corpus software, #LancsBox, which includes an innovative module for visualisation of collocations (Brezina et al., 2015; Liu, 2021). This approach is compared with the effects of learning collocations with a corpus-based collocations dictionary (Oxford Collocations Dictionary, McIntosh et al., 2009) and with non-corpus-based tools.

Both the effectiveness of a teaching method at improving language knowledge as well as the willingness of learners and teachers to use it are important for its successful implementation in language education. This study therefore sets out to examine both of these aspects – learners' linguistic gains together with their affective response after their experience with corpus-based learning. Furthermore, while there is growing evidence of the effectiveness of DDL in different contexts (e.g., Boulton & Cobb, 2017; Chen & Flowerdew, 2018; Lee et al., 2019), most of the research so far focused on the benefits of the method in the short term, by assessing learners' knowledge soon after they finished with using DDL. In addition to the short-term outcomes of DDL, this study will also investigate the longer-term effects of the method, by examining learners' ability to retain the knowledge gained through DDL as well as their plans for future use of corpus tools.

## 2. Literature review

## ***2.1 Role of collocations in L2 English***

The knowledge and appropriate use of collocations are regarded as a prerequisite of successful L2 production and a key component in L2 acquisition (Siyanova-Chanturia & Pellicer- Sánchez, 2019). First, collocations are very frequent across all modes of communication, accounting for nearly a third of both spoken and written discourse (Gablasova et al., 2017). Second, formulaic units, including collocations, play a major role in cognitive processes related to language use by enabling fluent processing and production of language in real time, with knowledge of collocations being linked to processing advantages in L1 as well as L2 (Myles & Cordier, 2017). Finally, the ability to use a range of collocations has a considerable impact on the complexity and appropriateness of learners' production in different linguistics domains such as in English for Academic Purposes (EAP) (Ackermann & Chen, 2013; Durrant, 2019).

Despite their importance in language use, collocations remain a very difficult component of language to acquire for L2 learners regardless of their level of proficiency (Boers et al., 2014; Gablasova et al., 2017; Men, 2018; Nesselhauf, 2004; Tsai, 2015). For example, as learners' proficiency increases, they tend to produce collocations more frequently. However, the range of these collocations remains limited, and non-target-like word combinations are still common (Granger, 2019). These issues can be attributed to the linguistic complexity of these multi-word units in terms of their semantic and syntactic properties, and the different degrees of fixedness they involve. For example, while some of the collocations (e.g., *make tea, take (a) shower, do homework*) are fairly fixed and act, to some extent, as a single lexical unit, others allow more flexibility and optionality in the combinations (e.g., *very/highly significant; I personally/strongly/firmly believe*). Transfer from learners' L1 can also contribute to non-target combinations of words (Webb & Kagimoto, 2011).

The second major factor, with an impact on the L2 acquisition and use of collocations, is related to the teaching approaches to these multi-word units. Collocations, together with other types of formulaic expressions, are often not given direct attention in language teaching and in textbooks, with only limited exposure provided to L2 learners compared with the teaching of single words (Pellicer-Sánchez & Boers, 2019; Tsai, 2015). Without enough exposure and guidance, L2 learners may find it difficult to develop an awareness of these formulaic units and, as a result, fail to either notice them in the input or create associative links between words (Wray, 2019). When collocations are given attention in teaching, traditional vocabulary exercises that rely on incidental learning (e.g., extensive reading), semi-incidental learning (e.g., input-enhancement techniques such as bolding or highlighting the target items) or intentional learning (e.g., decontextualized exercises such as matching of collocates) vary greatly in their effectiveness, frequently resulting in only marginal gains in linguistic knowledge (Pellicer-Sánchez & Boers, 2019).

## ***2.2 The effect of DDL on collocation learning***

Data-driven learning (DDL) is a corpus-based teaching approach that integrates data and findings from corpora into language teaching (Flowerdew, 2015). Corpora are electronic

collections of language samples, representing naturally-occurring language from different domains (e.g., academic writing) and from different groups of language users (e.g., L1 or L2 speakers, expert or novice writers in a domain) (McEnery & Hardie, 2011). In terms of theoretical support for the effectiveness of DDL, this teaching approach draws on two key principles: First, it uses corpora to identify patterns in target language use in order to expose students to multiple examples of the target patterns (e.g., grammatical structures, lexical regularities) embedded in their natural context and to draw their attention to them. This type of exposure enhances the likelihood that learners will notice the target patterns, which is a key prerequisite of learning (Schmidt, 2001). The noticing is further encouraged by the fact that corpus tools are designed to highlight the target pattern visually, for example, as part of placing the target feature in the centre of the concordance lines in the Key Word in Context (KWIC) display or as part of visualising the words that collocate (see Section 3.2 for examples). Different forms of input enhancement have been found to be crucial in guiding learners' attention to the target patterns (Smith, 1993). Given these characteristics, the corpus-based approach is especially valuable for teaching collocations, many of which are domain-specific and infrequent (e.g., *administer (a) test* or *unrequited love*), making it difficult for learners to encounter and notice them incidentally. Second, DDL represents an inductive, learner-centred learning approach based on the principles of constructivist theory of learning (Lee et al, 2019). In this type of learning, students are guided to construct their L2 knowledge themselves through interaction with language data (Collentine, 2000), with teachers acting as facilitators in order to develop learners' independence (Flowerdew, 2015). Specific DDL approaches can combine elements of discovery learning and inquiry-based learning, involving varying degrees of independent work by students, and guidance and scaffolding from the teachers (Charles, 2022). While discovery learning relies on minimal intervention from teachers and on independent work by students, inquiry-based learning involves more systematic involvement from the teachers throughout the learning process (Ozdem-Yilmaz & Bilican, 2020). Due to the level of cognitive processing involved in DDL, knowledge and skills gained by students tend to be more robust and retained for longer as learners have to actively engage with the learning process, rather than being presented with the target information (Lee et al, 2020; Collentine, 2000).

DDL has been found to lead to gains across a range of linguistic features and with different learner populations (Boulton & Cobb, 2017; Lee et al., 2019). Studies that looked at its use in collocation learning reported that students engaged in DDL usually outperformed learners taught using other methods, such as dictionary-based activities (Daskalovska, 2015) or textbook exercises (Rezaee et al., 2015; Vyatkina, 2016a). Using corpora as reference tools was also considered to show great potential, especially in identifying lexical errors (Chang, 2014; Bridle, 2019). In terms of corpus methods, the vast majority of these DDL studies relied on concordancing, that is, presenting multiple occurrences of the target collocations in their immediate context (e.g., Chan & Liou, 2005; Daskalovska, 2015; Liu & Ma, 2011; Rezaee et al., 2015; Vyatkina, 2016a), a very popular technique in DDL research and practice (Chen & Flowerdew, 2018).

In addition to the perceived effectiveness of a teaching method, a positive experience with the method is also likely to affect students' and teachers' willingness to adopt it on a regular basis (Chambers, 2019; Luo, 2016). So far, there have been mixed findings regarding learners' and teachers' attitudes towards using corpora in classrooms. While the novelty and

effectiveness of the method have been appreciated, negative aspects of DDL have also been pointed out, including the level of technical skills required, the time-consuming nature of corpus-based activities and the logistic difficulties with accessing computers during class-time (Chen & Flowerdew, 2018; Geluso & Yamaguchi, 2014; Lin & Lee, 2019; Luo, 2016; Yoon & Hirvela, 2004). Difficulties with interpreting the patterns from corpora also contributed to the negative experience of some students (Geluso & Yamaguchi, 2014; Yoon & Hirvela, 2004). For example, Yoon and Hirvela (2004) investigated L2 learners' attitudes towards their use of online Collins COBUILD corpus in a university writing classroom over a period of 10 weeks. The students found that the corpus was particularly useful for acquiring lexical patterns and enhancing their writing skills. However, they also reported difficulties with corpus consultation which was most frequently related to the time-consuming data analysis and understanding the cut-off concordance sentences. These issues continue to raise questions about the potential of corpus consultation to be successfully integrated into mainstream language pedagogy (Chambers, 2019).

### ***2.3 The current study***

This study seeks to extend our understanding of the effect of corpus-based teaching approaches on collocation learning by focusing on three key areas: First, it aims to capture both short-term and longer-term effects of the DDL on collocation learning, following the calls for the need to better understand learners' retention of knowledge once the DDL intervention has finished (Boulton & Cobb, 2017). While previous DDL research on collocation learning focused on receptive/productive knowledge of collocations, the current study also examines another dimension of collocation knowledge - learners' certainty about which words collocate. Learners' confidence has been shown to significantly affect their willingness to communicate and use their language knowledge productively (Lee & Lee, 2020); it is thus likely to be a major factor in their willingness to use the collocations they know in their writing and speaking. Second, the study explores students' experience with the use of DDL and corpus-based tools, both during and after the study. Although previous studies investigated students' attitudes towards DDL, little is known about learners' engagement with corpus tools beyond the intervention period (Charles, 2014). Learners' perception of DDL is an important factor when evaluating the potential of this method to be integrated into mainstream language classes (Chen & Flowerdew, 2018). Finally, the majority of previous DDL studies on collocation learning employed concordancing as the main corpus technique used by learners to identify and analyse collocational patterns. While concordance-based exercises and materials have proved very effective, there have been calls for other forms of corpus pattern visualisation to be explored in language pedagogy as well (Anthony, 2019; Vyatkina & Boulton, 2017). This study, therefore, evaluates the potential of #LancsBox and collocational graphs as a new technique to be incorporated into corpus-based language learning. To offer insights into these areas, the study investigates the effectiveness of DDL for learning collocations and compares it with the use of a corpus-based collocations dictionary and a control group. It is guided by the following three RQs:

RQ1: What is the effect of DDL on L2 learning of collocations?

- a. What are the gains in the knowledge of collocations after the treatment?
- b. Are the learning gains retained three months later?

RQ2: What is the effect of DDL on learners' confidence when selecting word combinations which collocate (i.e. which are frequent and natural combinations in English)?

RQ3: What are these learners' perceptions of the DDL approach, compared with the use of a collocations dictionary?

### 3. Methodology

#### 3.1 Participants

The participants in this study were 100 third-year English-major undergraduates from a university in the middle-eastern part of mainland China. The students were also taking part in pre-service teacher training. Their English language proficiency ranged from upper-intermediate to advanced based on the Test for English Majors (TEM) (Level 4)<sup>1</sup> (M = 67.97, SD = 8.58), with the range of scores roughly corresponding to B2-C1 bands of the Common European Framework of Reference (CEFR), and on the Dialang Test<sup>2</sup> (vocabulary part), with all but one participant scoring at B2 or higher on the CEFR (the remaining participant scored at B1). The majority (N=95) of the participants were female, reflecting the gender distribution in the university course. As established by the pre-study questionnaire, none of the participants used either corpus tools or the Oxford Collocations Dictionary before. The participants were randomly divided into three groups defined according to the tools used in their collocation learning (see Table 1): i) a corpus software #LancsBox (the DDL group), ii) Oxford Collocations Dictionary (the OCD group) or iii) other tools of their choice (e.g., online bilingual dictionaries, translation apps and online searching) (the control group).

**Insert Table 1 here**

#### 3.2 Corpus tools: #LancsBox and corpora

The tool used by the DDL group was #LancsBox (version 3)<sup>3</sup> (Brezina et al., 2015), a desktop software for the analysis of corpus data, which allows analysing both, corpora included in it (e.g., the British National Corpus 2014) as well as users' own datasets. In terms of its practical use in the classroom, #LancsBox needs an Internet connection to be downloaded; after this, learners are able to use it offline on their devices (e.g., tablets and laptops). #LancsBox was

---

<sup>1</sup> Test for English Majors (TEM) (Level 4) is a criterion-referenced English language test specifically targeted at university undergraduates majoring in English Language and Literature in China, with four levels of performance, namely *excellent* (score 80 or above), *good* (score 70-79), *pass* (score 60-69) and *fail* (score below 60).

<sup>2</sup> DIALANG is a free online diagnostic language test aligned with CEFR (Alderson & Huhta, 2005).

<sup>3</sup> #LancsBox was developed in the ESRC Centre for Corpus Approaches to Social Science (CASS) at Lancaster University (<http://corpora.lancs.ac.uk/lancsbox/>) and has continued to be updated since then. The third version of the software, newly released at the time of the data collection, was used in the study.

chosen for the study as it represents a novel corpus tool which can extend the range of corpus resources that are already used in language teaching (Liu, 2021). It contains different functionalities (modules), enabling different types of corpus analysis. Several of these modules were used in this study. The GraphColl (“graphical collocations” tool), illustrated in Figure 1, is an innovative corpus tool that identifies and visualises collocations in the corpus by showing the strength of the association between two words (indicated by the length of the line) and the frequency of the collocate (indicated by the colour density of the blue circle) (Brezina et al., 2015). Collocation graphs offer a clear visualisation of collocates related to the node word, thus directing learners’ attention to the relationship between pairs of words (Brezina, 2018; Liu, 2021). This feature enables learners to actively analyse the available linguistic data and identify collocational patterns (e.g., the strength of association between two words) by themselves. These activities make collocation learning with tools such as #LancsBox fundamentally different from consulting a dictionary, which presents learners with a list of target collocations and, in some cases, also example sentences illustrating their use.

**Insert Figure 1 here**

The second module of interest to this study, the KWIC (‘Key Word in Context’), allows conducting concordance analysis, in which all occurrences of the target linguistic feature (e.g., a word or a linguistic structure, such as the passive voice) are searched for and displayed in their immediate context (see Figure 2). Concordancing tools such as the KWIC module have been found to be very effective in corpus-based language teaching (Boulton & Cobb, 2017). The GraphColl module makes it very easy for learners to explore the context in which the collocations identified in the graphs naturally occur; users only need to right-click on a collocate to see it displayed in a set of concordance lines in the KWIC view. While concordancing tools are commonly used in corpus-based collocation learning, further visualisation of the links between the words is usually not available.

**Insert Figure 2 here**

Whilst the KWIC and GraphColl modules could be considered especially useful for L2 learners, additional functionalities in #LancsBox can further support language learning. For example, the split-screen-view (see Figure 3) allows comparing patterns from two corpora side-by-side and can be used to contrast learners’ own production with expert writing.

**Insert Figure 3 here**

In this study, two corpora representing written academic language were used by learners to look up and explore collocations in #LancsBox: the British Academic Written English (BAWE) corpus (Alsop & Nesi, 2009) and the British National Corpus (BNC) (BNC Consortium, 2007). The whole BAWE corpus, containing approximately 6.5M words, was used; from the BNC, the Arts sub-corpus was selected containing over 1M words and representing academic written language across a range of disciplines from the Humanities.

### 3.3 Oxford Collocations Dictionary

The computer-based Oxford Collocations Dictionary for Students of English (McIntosh et al., 2009) served as the second major tool in this study (Figure 4). The OCD was chosen for two reasons: First, it is a corpus-based reference work designed for upper-intermediate to advanced learners to assist with their writing and contains approximately 9,000 collocation headwords along with example sentences adapted from the Oxford English Corpus. Second, from a practical perspective, the dictionary offers easy access to an offline computer-based resource for checking collocations. While the dictionary was based on corpus evidence to identify and select collocations to be included, it involves a very different type of learner engagement from that of DDL. For example, McGee (2012), who compared the use of collocation dictionaries (including the OCD) and DDL concordancing activities, noted that when using the dictionaries, the students are exposed to pre-selected collocations, without the opportunity to discover which (types of) words collocate. He concluded that “this is not at all inductive” (McGee, 2012, p. 352) and further activities would need to be designed to increase learner engagement when using collocations dictionaries for learning.

**Insert Figure 4 here**

### 3.4 Linguistic gains: Collocation test

A 100-item test was designed to measure learners’ receptive knowledge of collocations, before and after the treatment. The test focused on three common types of collocations: verb-noun (e.g., *draw (a) conclusion*) (55 items), adverb-verb (*contribute significantly*) (25 items) and adverb-adjective (*entirely clear*) (20 items) collocations. The proportion of the three collocation types reflects that in the Academic Collocation List (Ackermann & Chen, 2013) from which the items were chosen. These types of collocations are interesting from a pedagogical perspective as they are frequent in language but particularly challenging for L2 learners because of their varying degrees of fixedness and the effect of cross-linguistic interaction (e.g. transfer from L1) (Nesselhauf, 2004; Web & Kagimoto, 2011; Pérez-Paredes & Sánchez-Tornel, 2014). Drawing on previous research on testing collocation knowledge (Boers et al., 2014; Gyllstad & Schmitt, 2019), the test followed a binary-choice format, asking participants to decide which of the two pairs of words collocate with each other, e.g., *execute (an) interview* or *conduct (an) interview*, *strongly agree* or *greatly agree*, and *largely clear* or *fairly clear*. The distractor items were selected using two sources: i) previous studies that reported common collocational mistakes by L2 learners (e.g., Lu, 2017; Men, 2018), ii) they were proposed by two experienced EAP teachers at the university where the data collection was conducted. Students received the test on three occasions (before, immediately after, and three months after the treatment), with the order of items randomised in each test to minimise the effects of order and guessing. In addition, the students were asked to indicate how confident they felt about each choice using a scale with five points: 100, 75, 50, 25 and 0 per cent.



### ***3.5 Learner experience with corpus-based tools: Questionnaires***

To gain further insights into students' experience with the corpus tools, the participants were asked about different aspects of their corpus-based learning in a questionnaire administered once they completed the treatment. Another questionnaire was administered three months after the end of the treatment. The two questionnaires were part of a larger study on collocation learning and academic writing practices of university students in China (the full questionnaires are available from Author 1, 2019). Only the data from the sections of the questionnaires which addressed students' experience with the two corpus tools (#LancsBox and the OCD) are reported in this study. This part consisted of thirteen closed questions in the post-treatment questionnaire and one question in the delayed questionnaire. The questionnaire items from the relevant part of the questionnaires are reported in Section 4.3 and in Appendix 1. In the post-treatment questionnaire, the participants were asked to indicate their (dis)agreement with the statements (items) on a six-point Likert scale, i.e. to say whether they (1) strongly agree, (2) agree, (3) somewhat agree, (4) somewhat disagree, (5) disagree, and (6) strongly disagree. The questionnaire given to the participants three months after the end of the treatment, asked about their use of the two corpus tools at that point – Whether they have continued to use the tools.

### ***3.6 Procedure***

The study followed a quasi-experimental research design. The writing classes in the study were carried out as additional, optional writing practice for students, who were invited to take part in these sessions, which took place after their regular EAP writing class. The treatment was conducted in the classroom across the period of eleven weeks. ***Before the treatment***, students in two experimental groups installed the relevant software, #LancsBox or the OCD, onto their own laptops. They also completed the first collocation test. ***During the treatment***, students took part in three training sessions, and eight writing sessions spread across eleven weeks. As no prior knowledge of corpus-based tools was assumed, two one-hour training sessions were provided for the DDL and the OCD groups at the start (Week 1 and 2 of the treatment). Previous studies identified sufficient training as one of the factors that affected the success of DDL outcomes (Boulton, 2009; Yoon & Hirvela, 2004), so another training session was provided in the middle of the treatment to consolidate learners' knowledge (Week 7). No training was provided for the control group, who were writing their pieces following the practice in their regular classes and, when revising, were free to choose any resources they wished other than #LancsBox and the OCD. In the eight weekly, 1.5-hour writing sessions (Weeks 3-6 and 8-11), participants in all three groups first wrote an essay on a given topic in class and in the following session revised the draft using either the assigned tools (the OCD and DDL groups) or tools of their own choice (the control group). The topics were adapted from academic writing task for TEM, a genre which the students were familiar with. Participants were allowed to leave as soon as they finished working on their pieces of writing. This process was repeated four times during the treatment (i.e. the participants produced four initial drafts and four revised drafts in the course of the study). The participants were encouraged to attend to the use of collocations during their writing and revising, and were told that their writing would be collected for analysis.

The researcher was present in the classroom during the sessions in which students revised their drafts to assist in case of technical problems with the tools but no further support with collocation-related activities was provided. The researcher also ensured that the students in the experimental groups only used their assigned corpus-based reference tools and that the control group did not use either of these two tools. *After the treatment*, the participants completed the collocation test and a post-treatment questionnaire. Three months later, the same academic collocation test and a delayed post-treatment questionnaire were administered.

### **3.7 Data analysis**

To answer RQ1 about the L2 learning and retention of collocation knowledge, the collocation tests were scored with each correct answer given one point. Before the analysis, the data was checked to ensure that the relevant assumptions (e.g., the data being normally distributed) for the ANOVA test were met. The scores at the three occasions were compared using a repeated-measures ANOVA. *P*-values (using threshold of 0.05) and effect sizes (using partial eta-squared) were calculated, with  $\eta_p^2 = 0.01$  considered a small,  $\eta_p^2 = 0.06$  a medium, and  $\eta_p^2 = 0.14$  considered a large effect size (Cohen, 1988). The changes in collocation knowledge were analysed based on the whole test (see Author 1, 2019, for the discussion of the test scores with respect to different types of collocations). To answer RQ2, an average percentage was calculated for learners' degree of confidence about which words collocate, and a repeated-measures ANOVA was used to compare the scores. Regarding RQ3, the data from the questionnaires were analysed in terms of the degree of learners' agreement with each statement. Percentages, means and standard deviations were calculated based on the six-point Likert scale responses for each questionnaire item.

## **4. Results**

### **4.1 RQ1: Learning gains in collocation knowledge**

As shown in Table 2 (and visualised in Figure 5), there were only small differences in the initial mean test scores in the three groups.

**Insert Table 2 here**

The results reported in Table 2 show that the three groups exhibited different trends in the development of the test scores at the three measurement points. The DDL group saw a moderate increase in their scores after the treatment and the score remained stable three months later. A one-way repeated measures ANOVA showed the main effect of the type of treatment (DDL) on the test score [ $F(2, 66) = 7.238, p = .001$ ] with a large effect size ( $\eta_p^2 = .180$ ). A post-hoc Bonferroni test revealed that the pre-test mean score differed significantly from both the post-test and the delayed post-test scores ( $p = .003$  and  $p = .011$  respectively); the difference between the post-test and the delayed post-test was not statistically significant.

### Insert Figure 5 here

As for the OCD group, although the test score showed a small increase in the post-test, as displayed in Table 2, it declined in the delayed post-test. The analysis showed no significant effect of the OCD on the test score [ $F(2, 62) = .632, p = .535$ ]. Somewhat surprisingly perhaps, the control group showed an incremental increase in the test score across the three measurement points, with the effect of the treatment being statistically significant with a large effect size [ $F(2, 66) = 4.77, p = .012, \eta_p^2 = .126$ ]. A post-hoc Bonferroni test revealed that the difference in scores between the pre-test and the delayed post-test was statistically significant ( $p = .021$ ).

When comparing the differences between the three groups, a one-way repeated measures ANOVA revealed the main effect of the treatment type [ $F(2, 194) = 7.546, p < .001$ ], with a medium-size effect ( $\eta_p^2 = .072$ ). Post-hoc comparisons using the Bonferroni test showed that there was a statistically significant difference between the DDL and the control group ( $p = .049$ ) but not between the DDL and OCD groups. No statistically significant differences were found between the OCD and the control groups.

#### **4.2 RQ2: Learners' confidence in collocation knowledge**

RQ2 examined changes in how confident learners felt when selecting the correct collocation in the collocation tests over the three measurement points. Table 3 shows descriptive statistics for the degree of students' confidence, expressed as a percentage, when selecting the correct word combination; results are presented separately for i) correct and ii) incorrect answers. As shown in the table, learners in all groups initially expressed a similar degree of certainty about the correctness of their choice, indicating more certainty in the cases when they chose correctly (approx. 68 -71 per cent) than incorrectly (approx. 62 per cent). While the certainty of the OCD and the control group increased only to a little extent (by 2-3 percentage points) after the treatment, the DDL group showed an increase by approximately 9 percentage points, for both correct and incorrect answers. After three months, the confidence about the choice of the correct collocation in the OCD and control groups returned largely to the initial level while the DDL group saw only a marginal decline (2-4 points). A repeated-measures ANOVA showed that the difference between the initial level of confidence of the DDL group on the one hand, and at later points on the other hand was statistically significant for both the correct [ $F(2, 66) = 10.769, p < .001, \eta_p^2 = .246$ ] and incorrect [ $F(2, 66) = 12.035, p < .001, \eta_p^2 = .267$ ] answers, with both effect sizes being large. In both cases, post-hoc Bonferroni tests indicated that the difference between the initial levels of reported confidence and those reported in the post-test and the delayed post-test were statistically significant (all  $p < .05$ ). However, the difference between the level of confidence recorded in the post-test and the delayed post-test was not statistically significant. No statistically significant differences were found for the OCD and control groups between any of the measurement points.

### Insert Table 3 here

### ***4.3 RQ3: Learner experience with corpus-based learning of collocations***

RQ3 explored different aspects of learner experience with the corpus-based resources in the OCD and DDL groups using questionnaire data. The findings are discussed with respect to i) the perceived effectiveness of the corpus tools in learning collocations, ii) practical experience with using the tools, and iii) the use of the tools beyond the study context.

#### *Perceived effectiveness of corpus-based learning*

Figure 6 summarises the trends in learners' perception of the usefulness of the two corpus-based tools, #LancsBox and OCD, for learning collocations. The figure shows the percentage of learners from each group who agreed with the statement, i.e. they chose 'agree', 'somewhat agree' and 'strongly agree' on the questionnaire. (The same convention is applied in Figures 7 and 8 below). In general, the majority of learners in both the DDL group and the OCD group regarded both the corpus tool and the corpus-based dictionary as useful in collocation learning, with somewhat more positive responses from the DDL group (Items 1 and 2). Further explorations of the experience of the DDL group also showed positive experience with individual modules in #LancsBox when revising writing, with 94.1 per cent agreeing they found KWIC useful and 97.1 per cent reporting the same for GraphColl. In terms of ease of use, 88.2 per cent reported they could identify collocations easily in the KWIC module, while all learners indicated this was true of GraphColl.

**Insert Figure 6 here**

#### *Practical experience when using the tools*

Figure 7 displays the trends in learners' experience with the use of the tools and the effectiveness of the training. It is somewhat surprising that learners in the DDL group were more likely to consider the assigned corpus tool (#LancsBox) easy to use and reported fewer problems compared with the OCD group (Item 3), given that a dictionary could be considered a very typical reference resource in language learning. This result could be related to the effectiveness of the #LancsBox training sessions, as more participants from the DDL group found the training sufficient compared to the OCD group (Item 4). Learners from both groups considered the hands-on experience helpful in getting to know the tools better (Item 5), although around half reported difficulties with running the tools on their own devices (Item 6).

**Insert Figure 7 here**

#### *Using the tools outside of the study*

To better understand the potential of the corpus-based tools to be adopted by learners and integrated into their usual learning habits, we explored both learners' intentions as well as their actual use of these tools outside of the classroom. As is apparent from Figure 8, a large number of learners from both groups employed the tools outside of the classroom during the treatment weeks (Item 7). Further, nearly all learners in the DDL group (97.1%) and over 80 per cent in the OCD group agreed that they would continue using the tools in their future learning. With

respect to their future language teaching, all learners in the OCD group reported planning to use the tool, while nearly all learners in the DDL group did so (Item 9). Looking at the actual use of the tools after the classroom part of the study finished, the majority of the DDL group reported still using #LancsBox three months later, which was not the case for the OCD group where none of the students continued to use the dictionary (Item 10).

**Insert Figure 8 here**

## **5. Discussion**

In terms of learners' collocation knowledge, RQ1 results revealed that all groups made gains between the pre-test and the post-test, with the DDL group showing larger improvement than the control group, even if the gains appear relatively modest (on average, three additional collocations). This is in line with prior findings that collocations take a long time to acquire regardless of the teaching method (Chen, 2019; Granger, 2019; Men, 2018). When interpreting the findings, it is important to stress that, unlike in the majority of DDL studies on collocation learning (e.g., Rezaee et al., 2015; Vyatkina, 2016a), learners in this study were not taught and tested on their knowledge of a pre-defined set of collocations. Instead, they focused on collocations of their choice in academic writing and were given a test with a large number of items from the same domain to observe a more general increase in their explicit receptive knowledge of collocations (e.g., Chan & Liou, 2005), with individual learning gains being likely even larger (e.g., Author 1, 2019).

In general, the results support the findings from the previous studies on collocation learning which reported that DDL appeared more effective, to some extent, than other methods, such as decontextualized exercises and explicit learning of collocating pairs of words (e.g., Daskalovska, 2015; Koosha & Jafarpour, 2006; Liu & Ma, 2011; Rahimi & Momeni, 2012; Rezaee et al., 2015; Vyatkina, 2016a). Importantly, while the collocation knowledge in the OCD group declined, the gains made by the DDL group were found to be stable three months after the end of the writing sessions, indicating that DDL, with sufficient amount of initial training, followed by independent work by students, may indeed lead to a more robust knowledge and longer retention. Notably, this finding differs from previous DDL research that showed the immediate gains to decline after a short period of time (Çelik, 2011; Chan & Liou, 2005; Daskalovska, 2015; Vyatkina, 2016b). It is possible that learners' ability to retain the new collocation knowledge was related to the length of their engagement with corpus consultation (11 weeks) in the current study, that was considerably longer than in the previous studies which in many cases included only a few learning sessions (Boulton & Cobb, 2017). However, while previous DDL meta-analyses (e.g., Boulton & Cobb, 2017; Lee et al, 2019) suggested that a longer duration of the intervention (over 10 weeks) may be associated with greater linguistic gains, the link was not straightforward and the relationship between the length of the treatment and the retention of knowledge has not been systematically examined yet. Further, it is possible that the learning and retention advantage in the DDL group could be related to both, the more general DDL features (e.g., exposure to multiple examples of the target feature) as well as the specific type of DDL activities in this study. These involved, for

example, the personal relevance of the collocations explored by learners, productive use of these collocations in written output and the visualisations of the collocations and their strength through the GraphColl module. These features have been related to greater learning outcomes and motivation of L2 learners both in DDL and SLA research (e.g., Daskalovska, 2015; Izumi, 2003). It was also interesting to find that learners in the control group made small gradual gains in their knowledge across the three measurement points; this is most likely due to being allowed to use their favourite and familiar tools during the intervention, which they likely continued to work with also after the intervention part of the study ended.

Highlighting a previously under-studied aspect of collocation knowledge, RQ2 investigated changes in learners' confidence about which words collocate. This often remains an 'invisible' aspect of linguistic knowledge, in spite of the fact that it can be crucial for learners' willingness to use specific collocations productively (Lee & Lee, 2020). While all three groups showed similar confidence in their choices initially, after the treatment, the DDL group reported a significant increase in their confidence which remained high after three months. This finding has two important pedagogical implications. As discussed earlier, collocation knowledge develops very slowly especially in terms of broadening the range of collocations that are used productively, with learners confining themselves to a relatively narrow set even as their proficiency increases (Granger, 2019). Providing learners with necessary training and tools can enhance learners' confidence in using a greater range of collocations in their written production. From a cognitive perspective, it is also possible that the greater retention of knowledge observed in the DDL group is also related to the increased certainty. However, while the confidence of the learners' in the DDL group when making the correct choice increased to nearly 80 per cent, their certainty when selecting the dispreferred pair of words also rose by 10 points to over 70 per cent. Whereas it is encouraging that DDL can increase learners' confidence, this particular finding also highlights some of the general issues that have been raised with respect to this teaching method. The major ones concern the fact that self-directed and independent learning as part of DDL may result in incorrect conclusions being reached and held by learners as the method makes it more difficult for teachers to detect and correct these misconceptions. However, it should be stressed that once teachers are made aware of these issues, it is possible to effectively address them (e.g., through providing learners with additional training, teacher guidance, checks on understanding and feedback on language use) (Gablasova & Bottini, 2022).

As for RQ3, the results show that learners were generally positive towards using corpus tools for collocation learning and writing revision, which is consistent with previous findings (e.g., Bridle, 2019; Geluso & Yamaguchi, 2014; Rezaee et al., 2015; Vyatkina, 2016b). In general, the satisfaction rate of learners in the DDL group was higher than that of learners who used the dictionary. This trend could be related to the perceived effectiveness as well as the experience with the use of the two corpus resources (e.g., Daskalovska, 2015; McGee, 2012) with the novelty involved in learning collocations with #LancsBox (and with collocation graphs) having a positive impact on students' affective response.

While previous research indicated the pedagogical potential of corpus consultation outside the classroom (Chen & Flowerdew, 2018; Crosthwaite et al., 2019; Lin & Lee, 2019), few of these studies explored whether learners continued using corpora after the intervention finished (e.g., Charles, 2014). The current study thus provides crucial evidence of learners'

willingness to keep using corpus-based tools for self-directed learning. It is interesting, however, to see the discrepancy between the intended and the actual use reported by students three months after the writing sessions. While both the OCD and DDL groups expressed their intention to continue using the tools, only the DDL participants did so. The drop in the use of the collocations dictionary can be perhaps explained by the competition with alternative dictionary-like reference tools available (e.g., online bilingual English-Chinese dictionaries and apps) and learners reverting to the ones they were already familiar with. The continual use of #LancsBox demonstrates the viability of DDL and corpus consultation in language learning in the context of Chinese higher education, even with alternative online and computer-based tools available to the students. These findings also offer further evidence that longitudinal DDL intervention indeed can foster a high degree of autonomy and encourage self-directed learning (Crosthwaite et al, 2019). This finding is especially important as Chen and Flowerdew (2018) stress that the success of DDL in academic writing depends on the users' ability to use corpora independently rather than just in teacher-led sessions, since "academic writers need to move beyond the classroom to work with corpora on their own and deal with the various texts they need to produce in the real world" (p. 357).

In addition to the insights provided by this study, it is also important to consider its possible limitations and the extent to which these could help identify directions for further research. First, in order to ensure high ecological validity, the outcomes of the study could be partly affected by variables that usually cannot be strictly controlled in a classroom quasi-experimental research design, such as individual learner characteristics (e.g., motivation and language aptitude) or the use of the corpus resources outside of the classroom. While difficult to control given the nature of the research (e.g., longitudinal intervention), the information about whether learners engaged with the tools outside of the classroom was collected through the questionnaires, providing important information about the viability of corpus consultation as a standard learning method in real educational contexts. Second, when interpreting the findings, it was not possible to separate the effect of different factors that may have led to the learning and retention advantage in the DDL group (e.g., the length or the novelty of the learning process). This is a common feature of current DDL research, which has so far been mostly concerned with demonstrating the effectiveness of DDL as compared to other teaching methods. However, in future research, it would be useful to design studies that allow contrasting different types of DDL provision, directly investigating the effect of key variables identified in several comprehensive reviews of previous DDL research (Boulton & Cobb, 2018; Chen & Flowerdew, 2017; Lee et al, 2019). In particular, as highlighted by this study, it would be especially valuable from a pedagogical perspective to further investigate the effect of the length of the treatment on the durability of acquired (collocation) knowledge.

## **6. Conclusion**

The study brought evidence of the positive impact of corpus consultation using a corpus tool, #LancsBox, on different aspects of L2 collocation knowledge. In terms of methodological implications, it highlighted the importance of considering different aspects of L2 (collocation) knowledge (e.g., learners' confidence, the long-term retention) in order to more fully

understand the outcomes of pedagogical intervention. The findings of this study also have several implications for language pedagogy. In terms of implications for pedagogical practice, the study demonstrated that DDL, supported by sufficient training, can provide an effective method for learning complex linguistic features such as collocations. However, many effective teaching methods may not be successfully implemented in practice in mainstream classrooms due to the logistic or technical issues that make their use more complicated than simply bringing a textbook or a dictionary into the classroom. This issue has been raised repeatedly for DDL as well (Chambers, 2019; Chen & Flowerdew, 2018; Luo, 2016). This study offers further evidence that DDL can indeed be successfully integrated into mainstream teaching, demonstrating the viability of the method for both, learning in the classroom context as well as for autonomous learning by students outside of the classrooms.



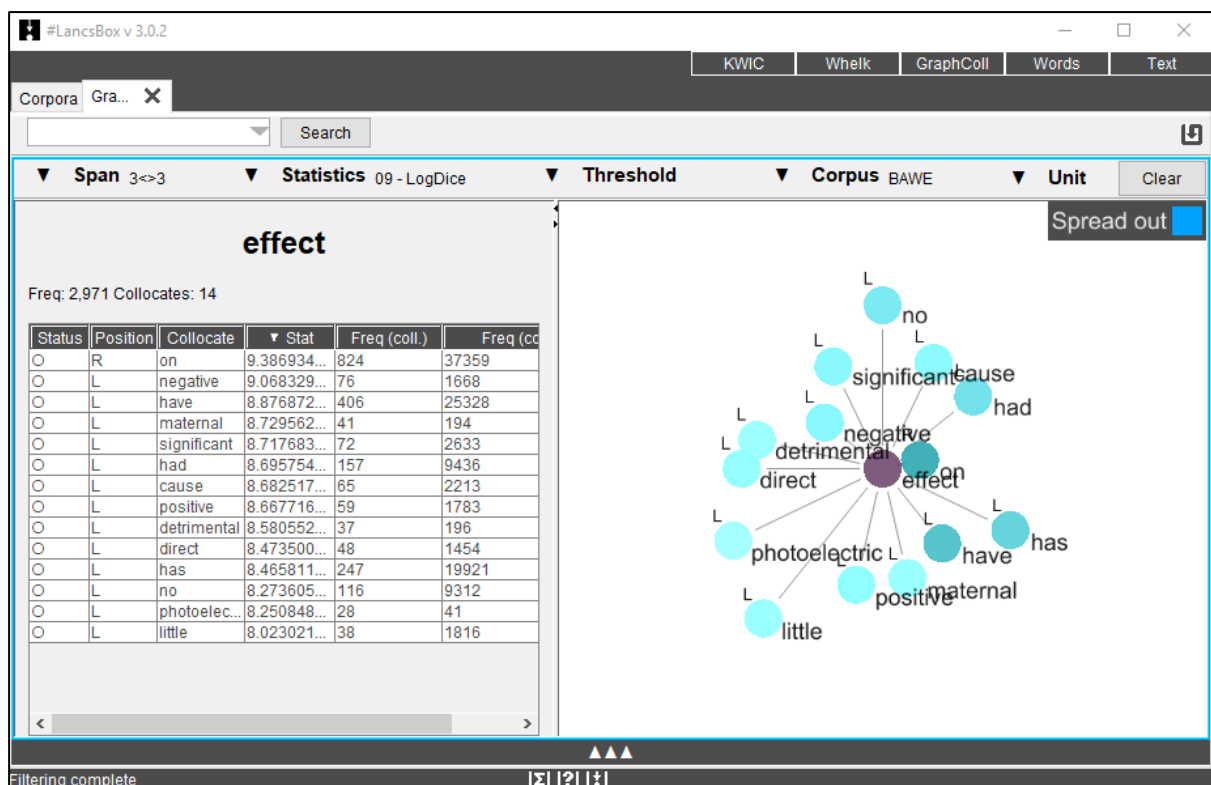
## References

- Ackermann, K., & Chen, Y. H. (2013). Developing the Academic Collocation List (ACL) - A corpus-driven and expert-judged approach. *Journal of English for Academic Purposes*, 12, 235–247.
- Alsop, S., & Nesi, H. (2009). Issues in the development of the British Academic Written English (BAWE) corpus. *Corpora*, 4(1), 71–83.
- Anthony, L. (2019). Tools and strategies for Data-Driven Learning (DDL) in the EAP writing classroom. In K. Hyland, & L. L.C. Wong (Eds.), *Specialised English: New directions in ESP and EAP research and practice* (pp. 179–194). Routledge.
- Boers, F., Demecheleer, M., Coxhead, A., & Webb, S. (2014). Gauging the effects of exercises on verb-noun collocations. *Language Teaching Research*, 18(1), 54–74.
- Boulton, A. (2009). Testing the limits of data-driven learning: Language proficiency and training. *ReCALL*, 21(1), 37–54.
- Boulton, A. & Cobb, T. (2017). Corpus use in language learning: A meta-analysis. *Language Learning*, 67(2), 348–393.
- Brezina, V. (2018). Collocation Graphs and Networks: Selected Applications. In P. Cantos-Gómez & M. Almela-Sánchez (Eds.), *Lexical collocation analysis: Advances and applications* (pp. 59–83). Springer.
- Brezina, V., McEnery, T., & Wattam, S. (2015). Collocations in context: A new perspective on collocation networks. *International Journal of Corpus Linguistics*, 20(2), 139–173.
- Bridle, M. (2019). Learner use of a corpus as a reference tool in error correction: Factors influencing consultation and success. *Journal of English for Academic Purposes*, 37, 52–69.
- British National Corpus, version 3 (BNC XML Edition). (2007). Distributed by Bodleian Libraries, University of Oxford, on behalf of the BNC Consortium. URL: <http://www.natcorp.ox.ac.uk/>
- Çelik, S. (2011). Developing collocational competence through web based concordance activities. *Novitas-ROYAL (Research on Youth and Language)*, 5(2), 273–286.
- Chambers, A. (2019). Towards the corpus revolution? Bridging the research–practice gap. *Language Teaching*, 52(4), 460–475.
- Charles, M. (2022). Corpora and autonomous language learning. In R. R. Jablonkai & E. Csomay (Eds.) *The Routledge Handbook of Corpora and English Language Teaching and Learning* (pp. 406–419). Routledge.
- Chang, J. Y. (2014). The use of general and specialized corpora as reference sources for academic English writing: A case study. *ReCALL*, 26(2), 243–259.
- Charles, M. (2014). Getting the corpus habit: EAP students' long-term use of personal corpora. *English for Specific Purposes*, 35, 30–40.
- Chan, T. P., & Liou, H. C. (2005). Effects of web-based concordancing instruction on EFL students' learning of verb-noun collocations. *Computer Assisted Language Learning*, 18, 231–251.
- Chen, W. (2019). Profiling collocations in EFL writing of Chinese tertiary learners. *RELC Journal*, 50(1), 53–70.

- Chen, M., & Flowerdew, J. (2018). A critical review of research and practice in data-driven learning (DDL) in the academic writing classroom. *International Journal of Corpus Linguistics*, 23(3), 335-369.
- Cohen, J. (1988). *Statistical Power Analysis for the Behavioral Sciences*. Routledge.
- Collentine, J. (2000). Insights into the construction of grammatical knowledge provided by user-behavior tracking technologies. *Language Learning & Technology*, 3(2), 44-57.
- Crosthwaite, P., Wong, L. L., & Cheung, J. (2019). Characterising postgraduate students' corpus query and usage patterns for disciplinary data-driven learning. *ReCALL*, 31(3), 255–275. <https://doi.org/10.1017/s0958344019000077>
- Daskalovska, N. (2015). Corpus-based versus traditional learning of collocations. *Computer Assisted Language Learning*, 28(2), 130–144.
- Durrant, P. (2019). Formulaic language for English for academic purposes. In A. Siyanova-Chanturia & A. Pellicer-Sánchez (Eds.), *Understanding formulaic language: A second language acquisition perspective* (pp. 174–191). Routledge.
- Flowerdew, L. (2015). Data-driven learning and language learning theories: Whither the twain shall meet. In A. Leńko-Szymańska & A. Boulton (Eds.), *Multiple affordances of language corpora for data-driven learning* (pp. 15–36). John Benjamins.
- Gablasova, D. & Bottini, R. (2022). Spoken Learner Corpora for Language Teaching. In Jablonkai, R. & Csomay, E. (Eds.) *The Routledge Handbook of Corpora and English Language Teaching and Learning*. Routledge.
- Gablasova, D., Brezina, V. & McEnery, T. (2017). Collocations in corpus-based language learning research: Identifying, comparing, and interpreting the evidence. *Language Learning*, 67, 155–179. doi:10.1111/lang.12225.
- Geluso, J., & Yamaguchi, A. (2014). Discovering formulaic language through data-driven learning: Student attitudes and efficacy. *ReCALL*, 26(2), 225–242.
- Granger, S. (2019). Formulaic sequences in learner corpora. In A. Siyanova-Chanturia & A. Pellicer-Sánchez (Eds.), *Understanding formulaic language: A second language acquisition perspective* (pp. 228-238). Oxon: Routledge.
- Gyllstad, H., & Schmitt, N. (2019). Testing formulaic language. In A. Siyanova-Chanturia & A. Pellicer-Sánchez (Eds.), *Understanding formulaic language: A second language acquisition perspective* (pp. 174-191). Routledge.
- Izumi, S. (2003). Comprehension and production processes in second language learning: In search of the psycholinguistic rationale of the output hypothesis. *Applied Linguistics*, 24(2), 168-196.
- Koosha, M., & Jafarpour, A. A. (2006). Data-driven learning and teaching collocation of prepositions: The case of Iranian EFL adult learners. *Asian EFL journal*, 8(4), 192–209.
- Lee, J. S., & Lee, K. (2020). Affective factors, virtual intercultural experiences, and L2 willingness to communicate in in-class, out-of-class, and digital settings. *Language Teaching Research*, 24(6), 813-833.
- Lee, H., Warschauer, M., & Lee, J. H. (2019). The effects of corpus use on second language vocabulary learning: A multilevel meta-analysis. *Applied Linguistics*, 40(5), 721-753.
- Lin, M. H., & Lee, J. Y. (2019). Pedagogical suitability of data-driven learning in EFL grammar classes: A case study of Taiwanese students. *Language Teaching Research*, 23(5), 541-561.

- Liu, S., & Ma, Z. (2011). An empirical study on concordance-based English collocation teaching. *I.J. Education and Management Engineering*, 4, 46–52.
- Liu, T. (2021). Data-driven learning: Using #LancsBox in academic collocation learning. In P. Pérez-Paredes & G. Mark (Eds.). *Beyond concordance lines: Applications of corpora in language education* (pp. 177–206). John Benjamins.
- Lu, Y. (2017). *A corpus study of collocation in Chinese learner English*. Routledge.
- Luo, Q. (2016). The effects of data-driven learning activities on EFL learners' writing development. *SpringerPlus*, 5(1), 1255.
- McEnery, T., & Hardie, A. (2011). *Corpus linguistics: Method, theory and practice*. Cambridge: Cambridge University Press.
- McIntosh, C., Francis, B., & Poole, R. (Eds.) (2009). *Oxford Collocations Dictionary for Students of English* (2nd ed.). Oxford University Press.
- McGee, I. (2012). Collocation Dictionaries as Inductive Learning Resources in Data-Driven Learning—An Analysis and Evaluation. *International Journal of Lexicography*, 25(3), 319-361.
- Men, H. (2018). *Vocabulary increase and collocation learning: A corpus-based cross-sectional study of Chinese learners of English*. Springer.
- McIntosh, C., Francis, B., & Poole, R. (Eds.) (2009). *Oxford Collocations Dictionary for Students of English* (2nd ed.). Oxford University Press.
- Myles, F., & Cordier, C. (2017). Formulaic sequence (FS) cannot be an umbrella term in SLA: Focusing on psycholinguistic FSs and their identification. *Studies in Second Language Acquisition*, 39(1), 3-28.
- Nesselhauf, N. (2004). Learner corpora and their potential for language teaching. In J. Sinclair (Ed.). *How to use corpora in language teaching* (pp. 125–152). John Benjamins.
- Ozdem-Yilmaz, Y., & Bilican, K. (2020). Discovery Learning—Jerome Bruner. In *Science Education in Theory and Practice* (pp. 177-190). Springer.
- Pérez-Paredes, P., & Sánchez-Tornel, M. (2014). Adverb use and language proficiency in young learners' writing. *International Journal of Corpus Linguistics*, 19(2), 178-200.
- Pellicer-Sánchez, A., & Boers, F. (2019). Pedagogical approaches to the teaching and learning of formulaic language. In Siyanova-Chanturia, A. & Pellicer-Sánchez, A. (Eds.), *Understanding Formulaic Language: A Second Language Acquisition Perspective* (pp. 153–173). Routledge.
- Rahimi, M. & Momeni, G. (2012). The effect of teaching collocations on English language proficiency. *Procedia-Social and Behavioral Sciences*, 31, 37–42.
- Rezaee, A. A., Marefat, H. & Saeedakhtar, A. (2015) Symmetrical and asymmetrical scaffolding of L2 collocations in the context of concordancing, *Computer-Assisted Language Learning*, 28(6), 532–549.
- Schmidt, R. (2001). Attention. In P. Robinson (Ed): *Cognition and Second Language Instruction* (pp. 3–32). Cambridge University Press.
- Siyanova-Chanturia, A. (2015). Collocation in beginner learner writing: A longitudinal study. *System*, 53, 148–160.
- Siyanova-Chanturia, A., & Pellicer- Sánchez, A. (Eds.). (2019). *Understanding formulaic language: A second language acquisition perspective*. Routledge.

- Smith, M. S. (1993). Input enhancement in instructed SLA: Theoretical bases. *Studies in second language acquisition*, 15(2), 165-179.
- Tsai, K. J. (2015). Profiling the collocation use in ELT textbooks and learner writing. *Language Teaching Research*, 19(6), 723–740.
- Vyatkina, N. (2016a). Data-driven learning for beginners: The case of German verb-preposition collocations. *ReCALL*, 28(2), 207–226.
- Vyatkina, N. (2016b). Data-driven learning of collocations: Learner performance, proficiency, and perceptions. *Language Learning and Technology*, 20(3), 159–179.
- Vyatkina, N., & Boulton, A. (2017). Corpora in language learning and teaching. *Language Learning and Technology*, 21(3), 1–8.
- Webb, S., & Kagimoto, E. (2011). Learning collocations: Do the number of collocates, position of the node word, and synonymy affect learning? *Applied Linguistics*, 32(3), 259–276.
- Wray, A. (2019). Concluding question: Why don't Second Language Learners More Proactively Target Formulaic Sequences? In Siyanova-Chanturia, A., & Pellicer-Sánchez, A. (Eds.) *Understanding formulaic language: A second language acquisition perspective* (pp. 248-269). Routledge.
- Yoon, H., & Hirvela, A. (2004). ESL student attitudes towards corpus use in L2 writing. *Journal of Second Language Writing*, 13(4), 257–283.



**Figure 1.** Interface of GraphColl in #LancsBox by using “effect” as an example

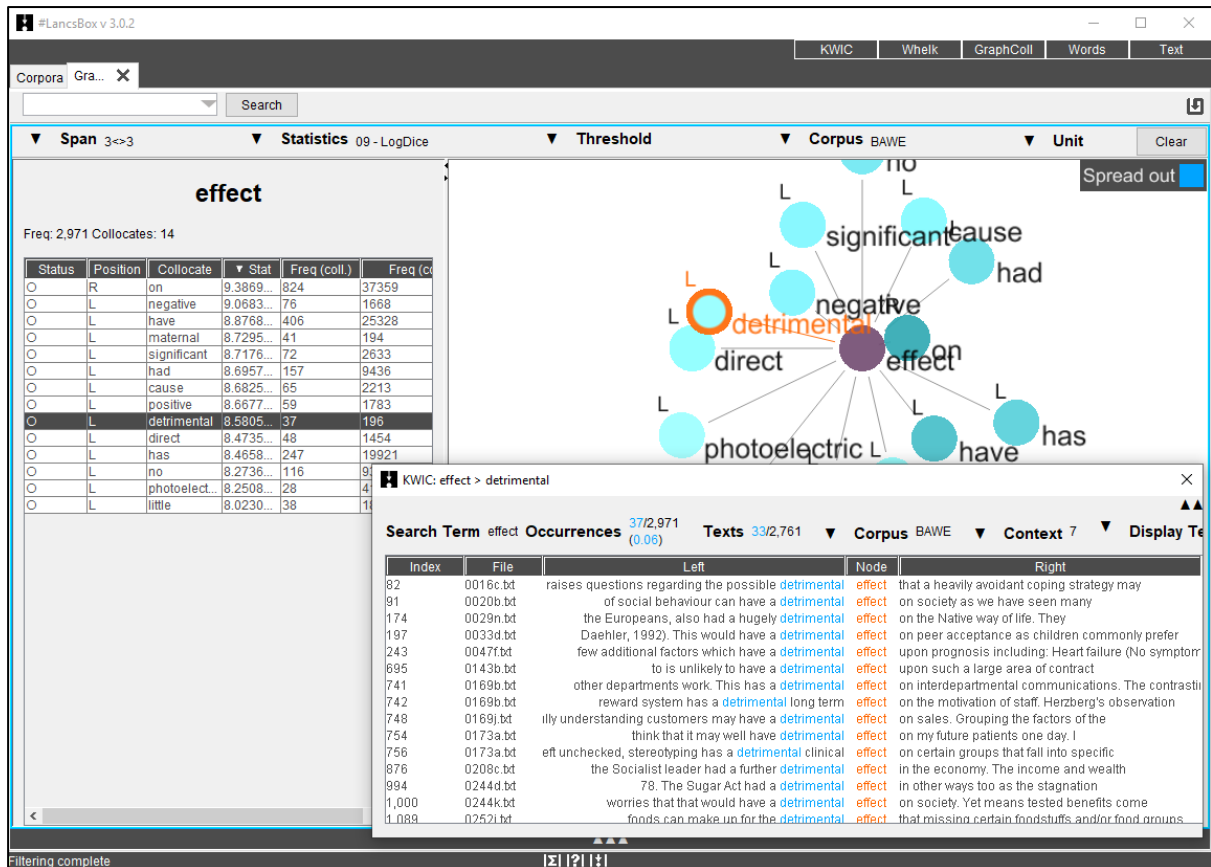


Figure 2. Interface of GraphColl linked with KWIC by using “effect” and “detrimental” as an example

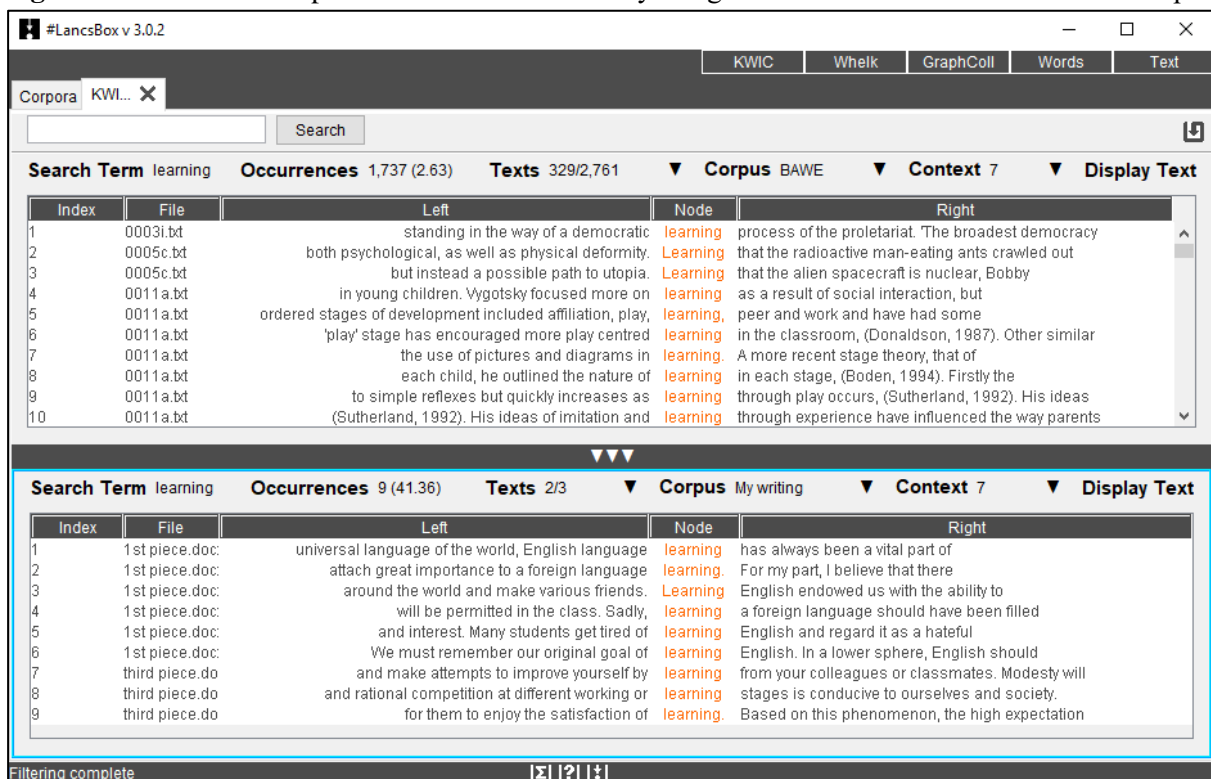


Figure 3. Interface of the split screen when using KWIC by searching “learning” as an example



Figure 4. Interface of the Oxford Collocations Dictionary using “responsibility” as an example

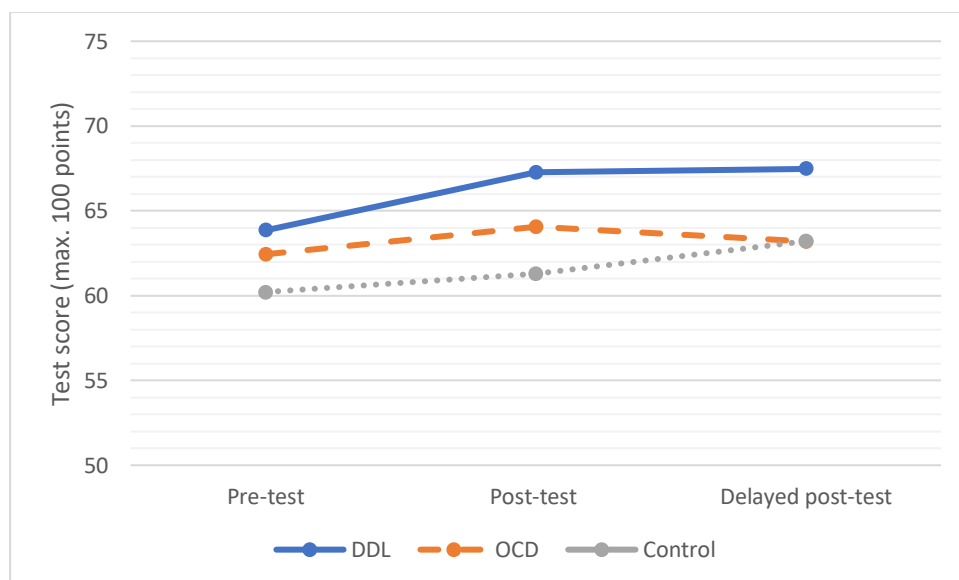
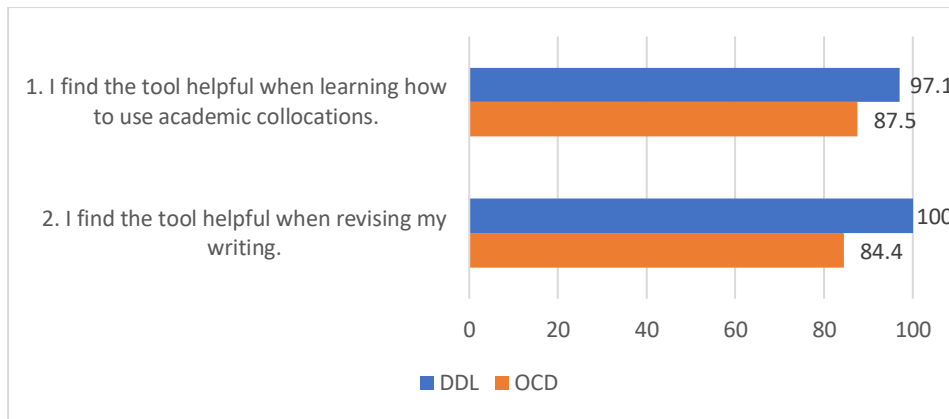
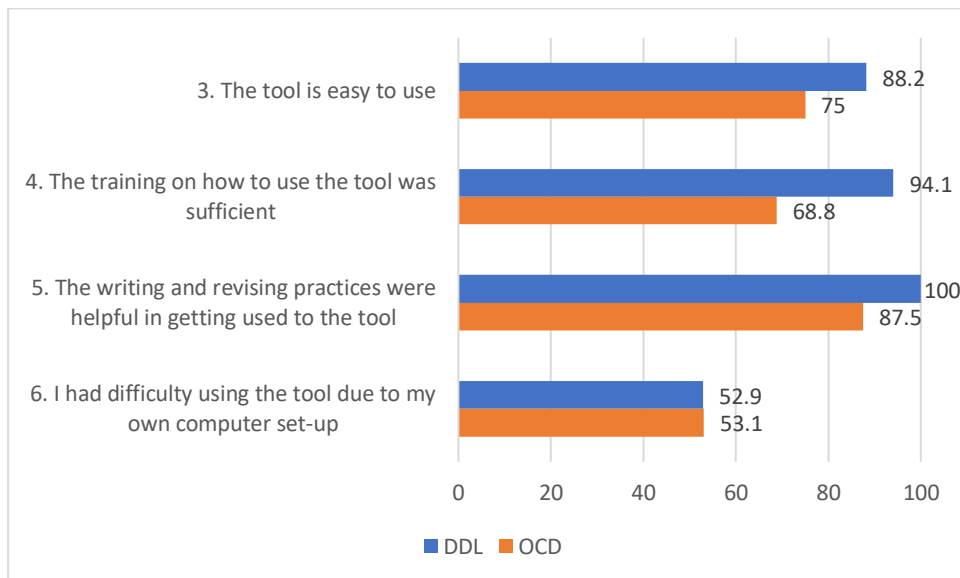


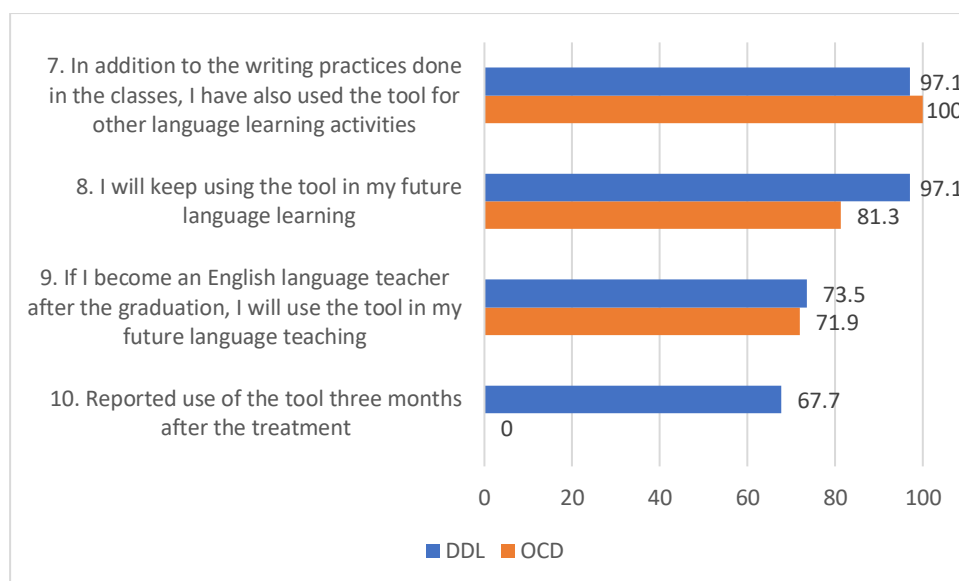
Figure 5. Mean score on the collocation test for the three groups



**Figure 6.** Learners' perception of the usefulness of DDL/OCD in learning collocations



**Figure 7.** Learners' attitudes towards the tools and the treatment



**Figure 8.** Learners' perceptions of the potential of the tools

**Table 1.** Information about participants in the three different groups

Group	N	Average age (SD)	Average years of English language learning (SD)	Teaching approaches (tools used)
DDL	34	20.26 (0.73)	11.23 (2.60)	#LancsBox (v. 3.0) (Brezina et al., 2015) and selected academic corpora
OCD	32	20.63 (0.92)	10.74 (1.93)	Computer-based OCD (McIntosh et al., 2009)
Control	34	20.65 (1.08)	10.44 (2.13)	Any available tools except for #LancsBox and the OCD

**Table 2.** Collocation test scores across the groups (mean and SD)

Group	N	Pre-test		Post-test		Delayed post-test		df	F	p	$\eta_p^2$
		M	SD	M	SD	M	SD				
DDL	34	63.85	7.83	67.26	8.78	67.47	10.64	2	7.238	.001	.180
OCD	32	62.44	6.63	64.06	9.21	63.19	10.98	2	.632	.535	.020
Control	34	60.21	7.52	61.29	8.07	63.21	7.97	2	4.77	.012	.126

**Table 3.** Confidence (mean percentage) reported when selecting collocations in the test

Group	Pre-test		Post-test		Delayed post-test		df	F	p	$\eta_p^2$
	M	SD	M	SD	M	SD				



Correct answer										
DDL	70.74	13.51	79.55	12.32	76.01	16.72	2	10.769	<.001	.246
OCD	68.21	16.86	71.64	15.52	67.87	15.54	2	1.459	.24	.045
Control	70.09	15.46	72.95	15.28	71.27	15.19	2	.700	.500	.021
Incorrect answer										
DDL	62.15	14.70	72.09	13.24	70.79	16.84	2	12.035	<.001	.267
OCD	62.36	18.04	64.20	17.52	62.62	16.02	2	.321	.727	.010
Control	63.14	16.16	64.47	16.44	63.91	16.24	2	.251	.779	.008