

Discovering formulaic language through data-driven learning: Student attitudes and efficacy

JOE GELUSO

*Iowa State University, USA
(email: joe.geluso@gmail.com)*

ATSUMI YAMAGUCHI

*Kanda University of International Studies, Japan
(email: yamaguchi-a@kanda.kuis.ac.jp)*

Abstract

Corpus linguistics has established that language is highly patterned. The use of patterned language has been linked to processing advantages with respect to listening and reading, which has implications for perceptions of fluency. The last twenty years has seen an increase in the integration of corpus-based language learning, or data-driven learning (DDL), as a supporting feature in teaching English as a foreign / second language (EFL/ESL). Most research has investigated student attitudes towards DDL as a tool to facilitate writing. Other studies, though notably fewer, have taken a quantitative perspective of the efficacy of DDL as a tool to facilitate the inductive learning of grammar rules. The purpose of this study is three-fold: (1) to present an EFL curriculum designed around DDL with the aim of improving spoken fluency; (2) to gauge how effective students were in employing newly discovered phrases in an appropriate manner; and (3) to investigate student attitudes toward such an approach to language learning. Student attitudes were investigated via a questionnaire and then triangulated through interviews and student logs. The findings indicate that students believe DDL to be a useful and effective tool in the classroom. However, students do note some difficulties related to DDL, such as encountering unfamiliar vocabulary and cut-off concordance lines. Finally, questions are raised as to the students' ability to embed learned phrases in a pragmatically appropriate way.

Keywords: corpus, data-driven learning, formulaic language, language learning.

1 Introduction

1.1 Grammar, lexis, and phraseology

John Sinclair, a pioneer in the field of corpus linguistics and phraseology, opens the eleventh chapter of his book, *Trust the text* (2004), with a telling example of just how elusive grammar can be. Reflecting on the letterhead of a European society founded to promote the topic of phraseology, Sinclair notes the ambiguity of grammatical correctness in the English version of the society's name: *The European Society of Phraseology*. Being a

European society, the letterhead also appears in German and French as, *Europäische Gesellschaft für Phraseologie* and *Société Européenne de Phraséologie*, respectively. Sinclair (2004: 177) ruminates:

Notice that the preposition used in the English version is *of*, and when I first encountered this I felt it was, if not ungrammatical, certainly uncomfortable. In French the preposition is *de* and in German *für*. The regular translation of *de* in English is indeed ‘of’, but of *für* it is ‘for’. I wondered, does:

1. European Society for Phraseology

sound any better [than European Society of Phraseology]? Yes, I think it does, but I have no idea why.

Examples of this nature illustrate the advantage of viewing language at the syntactic or phrasal level as well as the lexical level. This is precisely where a phraseological view of language can prove helpful to language learners, as it eschews the traditional lexis / grammar dichotomy view of language in favor of a more integrated one. Phraseology is predicated on Sinclair’s (1991: 110) idiom principle, which is encapsulated in the simple observation that “a language user has available to him or her a large number of semi-preconstructed phrases that constitute single choices, even though they might appear to be analyzable into segments”. When language is taken in segments, or multi-word units of meaning, much of its ambiguity dissolves.

1.2 Formulaicity and terminology

The phraseology, or formulaicity, of language can be illuminated by corpus linguistics. Römer (2009: 141) observes that, “if there is one major finding of modern (computer) corpus linguistic research over the past 40 years, it is probably that language is highly patterned”. Hoey (2009: 36) takes this idea a step further, expounding that “grammar is the system that one falls back onto when the collocational and other patterns are not used”. Clearly, both of these statements allude to the pervasiveness of formulaicity in language. Erman and Warren (2000) estimate that 52.3% and 58.6% for written and spoken language, respectively, is formulaic. Considering the cognitive processing advantages associated with using patterned language over novel utterances, it is not surprising to see such substantial quantities of formulaic language in everyday discourse (Conklin & Schmitt, 2008; Jiang & Nekrsova, 2007; Tremblay, Derwing, Libben & Westbury, 2011).

For formulaic language, Ellis (2012: 27) identifies three broad qualities to keep in mind: frequency, association, and native norms. Frequency and strength of association are two measures typically applied to formulaic language. Frequency is a self-explanatory term that simply refers to how often words co-occur. One measure of strength of association between co-occurring words is MI value, which Kennedy (2008: 23) defines as “the actual frequency of co-occurrence of two words with the predicted frequency of co-occurrence of the two words if each were randomly distributed in the corpus”. Hunston (2002) proposes that an MI value of three or higher indicates a relatively strong collocation.

Just as language is replete with patterns, so is the field of applied linguistics with terminology for patterned language (Wray, 2002). For the purposes of this study we will adopt Wray’s terminology of *formulaic sequence / language*. Wray (2002: 9) defines a formulaic

sequence (FS) as “a sequence, continuous or discontinuous, of words or other elements, which appears to be prefabricated: that is, stored and retrieved whole from memory at the time of use”. This is an often-cited definition of formulaic language, and with good reason. The inclusion of *continuous* or *discontinuous* affords the term a great deal of flexibility in accepting a wide range of multi-word units, “from formulaic phrase, to limited-scope slot-and-frame pattern, to fully productive schematic pattern” (Ellis, 2012: 18).

1.3 Corpora in foreign language learning curricula

There are two common pedagogical applications of corpora in second/foreign language teaching and learning: indirect and direct applications (Römer, 2011). Indirect applications include researchers and teachers consulting corpora to inform curriculum and materials development, and may lead to authentic examples of language for textbooks rather than invented examples. Direct applications of corpora in language teaching and learning, on the other hand, typically involve learners accessing a corpus directly. This is perhaps most commonly identified with data-driven learning (DDL), a term coined by Tim Johns (Johns, 1986). Johns (1991: 30) defines DDL as “the attempt to cut out the middleman as far as possible and to give the learner direct access to the data”. The idea behind DDL is that learners act as language detectives, or researchers, investigating authentic examples of the target language on their own. Boulton (2010a: 535) explains that “learners are not taught overt rules, but they explore corpora to detect patterns among multiple language samples”. Hunston (2002: 170) contends that DDL supports learning because “students are motivated to remember what they have worked to find out”.

DDL appears to be generally well received by learners and theoretically sound as a language-learning tool. Ellis (2002: 144) reminds us that cognitive linguistic theory postulates that “all linguistic units are abstracted from language use”. In usage-based theories of language learning, frequency is crucial for acquisition because “‘rules’ of language, at all levels of analysis... are structural regularities that emerge from learners’ lifetime analysis of the distributional characteristics of the language input” (Ellis, 2002: 144). Gries (2008) suggests that there is a strong affinity between corpus linguistics and cognitive linguistics as they both rely heavily on frequency, and Boulton (2009: 39) maintains that “DDL... exploits processes that humans have evolved to be naturally good at: exposure to data, detection of patterns, extrapolation to other cases.” While differing from naturalistic first-language acquisition which is largely unconscious, DDL can be argued to be firmly grounded in cognitive linguistic theory as learners analyze masses of input in a quest to become more familiar with structural regularities via inductive means.

Studies that have made quantitative comparisons of the efficacy of DDL with more traditional approaches to teaching suggest that DDL leads to results which are at least as good as, if not better than, other approaches (e.g., Boulton, 2009; Boulton, 2010b; Cobb & Boulton, forthcoming). For example, Frankenberg-Garcia (2012), with a group of EFL learners in Portugal, compared the efficacy of using dictionary definitions to corpus examples with respect to (1) learning the meaning of a target word, and (2) learning how to appropriately use a target word on a syntactic level. Two of Frankenberg-Garcia’s hypotheses in the study were that dictionary definitions would be more effective in the comprehension of novel words, while corpus examples would be more effective in learning the proper usage of familiar words. The findings supported both hypotheses (see also Frankenberg-Garcia, this volume).

Moving away from the experimental to the more qualitative, many studies investigating corpus-based learning have focused on student attitudes and beliefs toward the approach

and/or the processes involved (Boulton, 2009: 38), and most pertain to EFL or ESL learners' writing (Chambers & O'Sullivan, 2004; Yoon & Hirvela, 2004; Chambers, 2005; Yoon, 2008; Chen & Baker, 2010; Kennedy & Miceli, 2010). While these studies report largely positive findings related to outcomes of corpus-based learning and student attitudes toward DDL, a number of drawbacks consistently emerge as well, such as lack of confidence with respect to the grammaticality of corpus findings, the time-consuming nature of DDL, and the difficulty of interpreting the results of corpus investigations (Chambers, 2005; Chambers & O'Sullivan, 2004; Yoon & Hirvela, 2004). Complaints of this nature have led a number of scholars to recommend substantial introduction and training in how to use corpora properly (Kennedy & Miceli, 2010; Yoon & Hirvela, 2004). Bernardini (2004: 26), for instance, recommends starting students with convergent tasks, that is, tasks that guide learners to the same outcome. Once learners become familiar with the interface, they can then move on to more divergent, or independent tasks.

A number of studies have investigated student attitudes and beliefs about corpus-based learning as a way to improve their writing, and while some delve into student attitudes and beliefs about corpus-based learning with respect to speaking (Aguado-Jiménez, Pérez-Paredes & Sánchez, 2012; Pérez-Paredes & Cantos Gómez, 2004), there are notably fewer of them. The aim of this study is three-fold. The first is to outline a course that puts DDL at the center of the curriculum with the aim of increasing learners' repertoires of formulaic language and their ability to employ FSS in conversation; the second is to gauge how effective students are in employing their target phrases in a pragmatically appropriate manner; and the third is to investigate student attitudes toward this approach to language learning.

2 The study

2.1 Context

The course described below was a semester-long optional course open to third- and fourth-year students in the Department of International Communication (IC) at a private foreign language university in Japan. It met twice a week for 15 weeks, with each session lasting 90 minutes. Both rooms in which the class met were equipped with laptop computers, one with 30 and the other with 15, and a wi-fi Internet connection.

2.2 Population

The class consisted of 30 students, ranging in age from 20 to 22 years old, of whom 21 were female and 9 male. All were Japanese, and 29 of the 30 students agreed to participate in the study. In accordance with departmental policy, these students had taken the Test of English for International Communication (TOEIC) and had scores ranging from 540 to 860, and a mean of 736. This corresponds roughly to A2/B2 in the Common European Framework of Reference for Languages (Council of Europe, 2001), or intermediate to mid-advanced levels (Educational Testing Service 2013).

2.3 Course syllabus

The first three weeks of the course were used to explain course aims, do reading and discussion activities about DDL and inductive learning versus deductive learning, and train

students in using the corpus via handouts with convergent tasks. The Corpus of Contemporary American English (COCA) (Davies 2008-) was chosen for the course as it is a large, publicly available corpus consisting of 450 million words. Also, the instructor was from the United States and felt more comfortable commenting on instances of American usage as opposed to another variety of English.

Beginning in the fourth week, attention turned to three main components of the course: speaking journals (see section 2.3.1), student-led lessons, and a final project. Students used COCA to investigate and discover FSs they wished to use in their speaking journals and teach their peers in weekly student-led lessons. Finally, the students took part in a project that had them conduct a Behavioral Profile (BP) study of near-synonymous words and phrases (see section 2.3.3). Due to space limitations, the majority of attention in this article will be devoted to the speaking journals.

2.3.1 Speaking journals. The speaking journals formed the core of the course, and students were responsible for completing four throughout the semester. Students were given four class periods over the course of two weeks to complete each speaking journal. The speaking journals consisted of four distinct phases: (1) preparation; (2) corpus consultation; (3) a rehearsal conversation; and (4) the real conversation.

The speaking journal task is essentially based on input and interaction. An interactionist perspective on language acquisition posits that “the interactional ‘work’ that occurs when a learner and his/her interlocutor (whether a native speaker or more proficient learner) encounter some kind of communication breakdown is beneficial for L2 development” (Mackey, Abbuhl & Gass, 2012: 9). In the course reported on here, the learners’ task was to identify a potential communication breakdown *before* it happened by learning FSs that they did not have command of prior to the corpus consultation. The learners then took their FSs and used them in conversation with a more proficient speaking partner. This use of novel FSs can reasonably be likened to Swain’s (1995) output hypothesis as the learners ‘pushed themselves’ to create the opportunity to use their target phrases. Productive use of the target language, Swain (1995) contends, causes learners to process the language more deeply than input alone.

Phase 1 of the speaking journal had the students interact with authentic materials that they were free to choose. It was hoped that choosing materials and topics that interested them would increase motivation. Furthermore, formulaic language is more ubiquitous in authentic materials, such as television and movies, than in textbooks designed for language learning (Irujo, 1986: 237; Biber, Conrad, and Cortes, 2004: 379–380). Students were given a number of choices for authentic materials: English-language video news (e.g., CNN news video), news articles (e.g., an online news source, or print newspaper), magazines, TV shows, movies, and comic books, or they could bring their own ideas to the instructor for approval. The students had easy access to the first two options through the Internet, and easy access to the remaining options through the university’s library and self-access center. They were not allowed to choose any one form of materials more than twice in order to guarantee exposure to a wide variety of media. After the students had read or watched their material, their task was to write up a summary and two discussion questions. The summary and discussion questions were then used for small-group discussions in the subsequent class period.

Phase 2 of the speaking journal was the DDL component, with students using COCA to investigate words and phrases discovered in Phase 1. A useful analogy to describe the goal

of this phase is Kennedy and Miceli's (2010: 32) *pattern-hunting*, which "amounts to encouraging them [the learners] to use the corpus as an aid to the imagination and memory". Learners were encouraged to investigate words that they anticipated would be of use in their planned topic of conversation. One student, for example, planned to discuss her search for a job as she was approaching the end of her university career. In her corpus investigations she began with the word *work*, as it was sure to come up in conversation. Working her way through the concordancing phase of the speaking journal she settled on the five-word phrase, *work on a full-time basis*. The student was then able to use this phrase in the 'real conversation' phase of her speaking journal. As it turned out, students in the class often chose to look up familiar words with the goal of finding novel ways to use them (recall Frankenberg-Garcia's (2012) study; see also Frankenberg-Garcia this volume). Upon choosing intriguing collocates, students noted the frequency and MI value, and combed through the concordance lines to find interesting patterns. If a student chose to investigate a phrase in COCA rather than an individual word, the collocate function of the corpus was not used. Figure 1 shows notes a student took in her speaking journal during her concordancing.

In Phase 3 of the speaking journal, students used the phrases from their previous investigations in small-group 'rehearsal conversations' with their classmates. The main point of these rehearsals was to give students the opportunity to practice using their new phrases in the context of their chosen topic, and, perhaps more importantly, to learn how to manipulate a conversation in order to create an opportunity to use their target phrases.

Phase 4 was the final phase of the speaking journal where students were sent out to record their 'real conversations' with a native or more proficient speaker of English. Students were afforded opportunities to practice conversation in an open space at the university where international students often gather and teachers are available for informal conversations. There are mp3 recorders available for students to borrow to record their conversations, but most students simply used a personal smartphone. Students began their real conversations with a topic and general plan as to how they anticipated working their FSs into their conversations based on their rehearsal. After the conversations, students completed a reflection section in their speaking journal where they listened to their recording as a whole, noted when they used their FSs, and whether or not they were able to produce the FSs as planned. The sound files of their recordings were then uploaded to the class website, or emailed to the instructor, and the speaking journals were handed in.

2.3.2 Student-led lessons. The second major component of the class was the student-led lessons. Starting in the fourth week of class, and each week thereafter, a small group of students led the class in a 30-minute lesson featuring FSs discovered through their speaking journals. Ten groups of three were formed, and each student contributed two of their favorite FSs to the lesson, so each lesson featured six FSs. In addition to explaining their FSs, students led their classmates in an activity designed to give the class an opportunity to use the FSs. Example activities include variations of Pictionary-like games where the class draws pictures of the target phrases for the other students to guess, telephone-like games, hot-potato, creating skits or writing stories that use the phrases, and so on. The student-led lessons proved to be very popular, as illustrated in the questionnaire results presented in

Speaking English Naturally

Original word 1

disaster

Interesting Collocates	Frequency	MI Value
the ^{edge (scary place)} brink of disaster	72	6.74
you can image easily the situation when you read or An unmitigated disaster listen.	53	9.64
If a situation spells disaster, it makes you expect spell (v) disaster disaster	67	5.74
humanitarian disaster	99	6.20

(marked by humanistic values.)

1. Good management practices in one could spell disaster in the other.
2. When more wrong stuffs arrive, it will spell disaster.
3. The view of town after hurricane is an unmitigated disaster.

Original word 2

utmost

Interesting Collocates	Frequency	MI Value
matter of ^(utmost) utmost ^{big} matter	20	3.29
I shall do my utmost ^{- try to do something the best}	6	4.00
really respect → have the utmost respect for	125	7.24
with the utmost care	60	4.25

take long time to think 1 thing carefully.

1. I need to speak to you of a matter of utmost importance.
2. I shall do my utmost to honor what best fits your schedule.
3. You should choose your partner with the utmost care.

Fig. 1. Example of a student’s notes from concordancing session on COCA

section 4.3. Some examples of FSs that students used in their speaking journals and then went on to teach their peers in class are:

- (1) pave the way for
- (2) behind the scenes
- (3) put (personal pronoun) best foot forward
- (4) poised on the brink of
- (5) catch one’s eye
- (6) fail to recognize
- (7) place an emphasis on

Table 1 *Examples of near-synonymous words and phrases investigated by students*

	Word/phrase 1	Word/phrase 2	Word/phrase 3
Example 1	huge	enormous	immense
Example 2	think	consider	wonder
Example 3	resemble	look like	similar
Example 4	expect to	hope for	look forward to
Example 5	look at	watch	see

2.3.3 Behavioral Profile study. The course culminated in students undertaking a Behavioral Profile (BP) study of near-synonymous words or phrases of their own choosing. Gries (2010) explains that BP studies allow for the fine-grained analysis of near synonyms, which can shed light on differences between near synonyms and polysemous words. The scope of the project was such that it is not possible to explore the details fully here, but briefly the project entailed students identifying near-synonymous words or phrases and embarking on a corpus analysis via COCA and web searches to reveal subtle differences in patterns of usage. Students wrote reports to present to classmates and finally hand in to the instructor. Table 1 provides examples of the type of near-synonymous words and phrases students investigated.

3 Data collection and analysis

In order to arrive at a clearer understanding about student attitudes toward this particular approach to corpus-based language learning, data were collected via a questionnaire and triangulated through follow-up interviews and student reflection logs at the end of each speaking journal.

The questionnaire consisted of 44 statements to which the respondents were asked to indicate their degree of agreement on a 6-point Likert scale. The majority of items were adapted from two published studies on using corpora in L2 writing (Yoon & Hirvela 2004; Liu & Jiang 2009). The researchers designed the remaining items specifically for this study. All statements were presented in English and Japanese. The questionnaire included negatively worded items to keep respondents from marking only one side of the questionnaire, and scores for such items were reverse-coded before analysis (Dörnyei & Taguchi 2010). The researchers merged items into multi-item scales based on theoretical considerations. Categories were: (1) difficulty in using corpora; (2) positive impact of using corpora; (3) effectiveness of presentation and delivery of coursework; (4) completing speaking journals and incorporating phrases; and (5) attitudes and beliefs about data-driven learning and its potential. Six of the participants did not answer all 44 items on the questionnaire, missing an item either by choice or simple oversight (e.g., participant 11 did not answer item 14; participant 1 did not answer item 23). Therefore, the internal consistency, or reliability in how participants responded between items on the questionnaire, was based on the responses of the 23 individuals who responded to all 44 items and measured via Cronbach's Alpha in the statistical package SPSS. The instrument showed a high level of internal consistency with $r = .870$. Means, modes, and standard deviations were calculated based on all participants' responses. For the purposes of presentation, the results from the questionnaire are presented simply in terms of agreement to the statements.

Table 2 *Positive impact of using corpus*

Item Category	N	Agree*	Mean**	Mode**	SD
18. Improved knowledge of collocations	29	28 (97%)	5.00	5	1.04
13. Helpful for learning the usage of vocabulary items	29	24 (83%)	4.62	5	1.21
14. Helpful for learning the usage of phrases	28	25 (89%)	4.68	5	1.09
15. Helpful for learning grammar	29	21 (72%)	3.97	4	1.21
7. Learned new phrases through familiar vocabulary	29	28 (97%)	4.86	5	1.03
26. Helpful to find new ways to use familiar vocabulary	29	27 (93%)	4.66	4, 5***	1.17
19. Learned new vocabulary	29	27 (93%)	5.00	5	1.10
41. Good for understanding the differences between near synonymous phrases	29	24 (83%)	4.45	5	1.12
12. More helpful than a dictionary for finding common phrases	29	20 (69%)	3.97	4	1.12
22. Helpful for writing	29	19 (66%)	3.79	4	1.01
23. Helpful for speaking	28	25 (89%)	4.36	4	0.99
21. Equally helpful for both speaking and writing	29	20 (69%)	3.86	4	1.09
20. Improved my English	28	21 (75%)	3.89	4	1.03
3. Helpful for language learning	29	28 (97%)	5.00	5	1.04

*Raw numbers and percentage in parentheses

**1: strongly disagree, 2: disagree, 3: somewhat disagree, 4: somewhat agree, 5: agree, 6: strongly agree

***Responses occurred equally.

The follow-up interviews were semi-structured with lead questions based on the survey results and student reflection logs from the speaking journals (see the Appendix for interview questions). Interviewees were selected at random and included two males and three females. Each interview lasted about thirty minutes. Interviewees were given the choice of being interviewed in English or Japanese, and all of them chose to communicate in Japanese. The quotes presented in this paper were translated into English by the researchers.

Finally, one additional quantitative analysis was performed to investigate whether students were able to employ their target phrases in a contextually and/or pragmatically appropriate manner. The analysis entailed giving a sample of 114 phrases to four native-speakers of English to independently rate on a numerical rating scale of 1–4 (1 being ‘inappropriate’ and 4 being ‘appropriate’).

4 Results and discussion

4.1 *Positive impact of corpus use*

The findings concerning the impact of corpus use were quite encouraging and in general suggest students’ belief in the utility of DDL on a number of fronts, as Table 2 demonstrates. Students felt strongly that this approach to language learning increased their knowledge of collocations. Nearly all participants agreed with the statements that researching familiar vocabulary items in the corpus led to learning new phrases and new ways to use familiar vocabulary (see items 7 and 26). This provides qualitative support to Frankenberg-Garcia’s (2012) finding that concordances are useful for learning novel usages of familiar words. Perhaps most encouraging is that 28 of 29 participants, 97%, believed DDL to be helpful

Table 3 *Difficulty with using corpus*

Item Category*	N	Difficult**	Mean***	Mode***	SD
2. Learning to use COCA	29	17 (59%)	3.66	4, 5****	1.29
1. Concordancing	29	14 (48%)	3.62	3	1.08
6. Finding phrases around key words	29	15 (52%)	3.69	3, 4, 5****	1.29
8. Unfamiliar vocabulary	29	22 (76%)	4.10	4	0.82
9. Cut-off sentences	29	24 (83%)	4.41	5	0.91
10. Too many sentences	29	13 (45%)	3.34	3, 4****	1.20
11. Limited access to internet	29	6 (21%)	2.83	3	1.39
27. Understand context of concordance lines	28	15 (54%)	3.79	3	1.07

*Agreement that the category is *difficult* as opposed to *easy*

**Raw numbers and percentage in parentheses

***1: strongly disagree, 2: disagree, 3: somewhat disagree, 4: somewhat agree, 5: agree, 6: strongly agree

****Responses occurred equally.

for language learning, with a mean score of 5.00. One student noted in her speaking journal log:

I had a good conversation with Shelley [a teacher]. It went as I planned. And I could learn new words from the conversation. I think it's one of the best learning styles. I think she [Shelley] has different ideas from rehearsal conversation.

Students also believed that COCA was helpful for writing (66%) and speaking (89%), (see items 22 and 23). The latter is particularly of note because, as mentioned earlier, less has been done in the way of investigating student beliefs about the benefits of corpus consultation as pertaining to speaking as compared to writing. This finding is perhaps not too surprising, though, as the course focus was on speaking rather than writing. Nevertheless, the students did perceive corpus consultation to be a useful tool to improve their speaking.

4.2 *Difficulty in using the corpus*

A recurring line in the literature is the difficulties that accompany DDL, and many common themes from previous studies emerged here (Table 3). One notable exception, and likely a sign of the times and location, is that very few participants believed a dearth of Internet access to be a hindrance. As noted earlier, a frequent complaint is the investment in learning how to use a corpus effectively. In this study, too, a slight majority of students felt that learning to use COCA was difficult. However, once past the initial learning curve, less than half of the students felt that the actual concordancing was "difficult". Interestingly, the participants were more or less split over the categorization of abundant concordance lines as a "difficulty", which traditionally has been a common complaint about corpus consultation. One interviewee, however, did explicitly note a common complaint with DDL, citing difficulties he had in understanding the meaning of concordances due to cut off sentences:

[In the concordance output] there are many example sentences. With the long sentences, I mean, I can't see the entire sentence, just one part... so I can't understand the "situation". If sentences are cut off in the middle... I wonder what the following words will look like...

Table 4 *Effectiveness of presentation and delivery of coursework*

Item Category	N	Agree*	Mean**	Mode**	SD
4. The training on how to use the corpus was necessary	29	26 (90%)	5.10	6	1.08
16. 90 minutes was sufficient for concordancing	29	14 (48%)	3.72	3	1.33
17. I believe I spent enough time concordancing for each speaking journal	29	24 (83%)	4.28	5	1.16
30. The rehearsal conversations were helpful	29	24 (83%)	4.72	5	1.30
31. The rehearsal conversations helped me to understand the context in which to use the key words	29	25 (86%)	4.41	4	1.18
39. The student-led lessons were an effective way to learn new phrases	29	27 (93%)	4.79	5	1.11
40. The student-led lessons were a fun way to learn new phrases	29	25 (86%)	4.72	5	1.28
34. Having a teacher or native speaker explain what I find in the corpus was helpful	29	28 (97%)	4.79	5, 6***	1.08
42. This class allowed me to direct my own learning	29	19 (66%)	3.76	4	0.99
37. This class afforded me many learning opportunities	29	26 (90%)	4.86	5	1.13

*Raw numbers and percentage in parentheses

**1: strongly disagree, 2: disagree, 3: somewhat disagree, 4: somewhat agree, 5: agree,

6: strongly agree

***Responses occurred equally.

In addition to cut-off sentences, a mismatch of register was also pointed out as a difficulty by an interviewee. This is of paramount importance, and likely one of the greatest weaknesses of the course reported on here. Students were using COCA to find high-frequency FSs to incorporate into their conversations, which usually happens in semi-informal contexts. Yet much of the spoken language accumulated in COCA comes from formal contexts, such as news programs. This underscores the importance of raising student awareness of the genre register from which concordance lines are gleaned and making decisions about how pragmatically appropriate a phrase from a news broadcast might be in a different context. This is a topic we will come back to in section 4.4.

4.3 *Effectiveness of presentation and delivery of coursework*

Given the unique nature of this class, the researchers wanted to gather data on the students' attitudes and beliefs about the delivery of the coursework. For this reason, a number of statements specifically addressing aspects unique to the context were crafted, such as items concerning class time allotted to concordancing, the rehearsal conversations, and student-led lessons (Table 4).

Because of the protracted nature of DDL, substantial time in class was allotted for concordancing. Typically students were given 60 minutes to concordance and then took part in a student-led lesson for the remaining 30 minutes. Occasionally, though, an entire 90-minute class period was devoted to concordancing and consulting with the teacher and their classmates about their findings. Even when given 90 minutes of class time, the majority of students felt that it was not enough. Interestingly, though, 83% believed that they did complete an adequate amount of concordancing for each speaking journal (items 16 and 17). Hopefully

this indicates that students spent time concordancing outside of class, but it may be that while 90 minutes was not enough to achieve their concordancing goals for a speaking journal, students felt that additional concordancing would not have been beneficial.

The rehearsal conversations were viewed in an overwhelmingly positive light. This is apparent from the questionnaire items related to the ‘helpfulness’ of the rehearsal. 24 of the 29 participants agreed that the rehearsal was helpful to some degree, with a mean score of 4.72 and a mode of 5. Likewise, a common theme in students’ speaking journal logs was the usefulness of the rehearsal in helping hone their conversations in order to use their planned FSs. One student wrote in her speaking journal log:

When I did a rehearsal conversation I couldn’t use the phrases well and I felt some phrases were unnatural. Therefore, I made more examples to be able to choose and use the most natural one while doing this conversation. I expect my [real] conversation to go more naturally than my rehearsal conversation.

Having a teacher or native speaker of English explain what students found in the corpus was also perceived positively, as illustrated by item 34. Two students reported in the follow-up interviews that they felt strongly about the need to check the meaning and usage of target FSs with a native speaker, teacher, or friend whose English was more advanced, because the meanings of new words and phrases encountered in the corpus were sometimes not found in electronic dictionaries, or the nuance of the target words and phrases was lost when translated into Japanese.

The student-led lessons also proved to be a popular activity throughout the course (see items 39 and 40 in Table 4). Indeed, the instructor of the class noted that students responded well to peers in the role of teacher, and that the students made substantial efforts to create engaging lessons. It is worth noting that the student-led lesson was weighted at 20% of the final grade, which may in part explain the effort the students put in.

4.4 Completing speaking journals and incorporating phrases

Wray and Fitzpatrick (2010: 38) point out that “it would be easy to construe them [FSs] as a straitjacket for the user, rather than an opportunity”. However, when used correctly in the appropriate context, FSs can be a concise, economical, ‘native-like’ means of conveying one’s message. Indeed, the more adept user of FSs can cut short, rearrange, and come up with new combinations joined by individual lexical items or shorter phrases. Mastery of a large repertoire of FSs can thus be seen as an integral step in the journey to fluency in a language. Wray and Fitzpatrick (2010) and Kennedy and Miceli (2010) suggest that considering the context and anticipating the trajectory of a planned conversation when choosing target phrases (i.e., pattern-hunting) will ostensibly minimize the failure to employ pre-determined phrases in a conversation. The following student comment illustrates her success with this approach:

I visualized and simulated the trajectory of a conversation using the target FSs. In addition, I tried to visualize how to expand the conversation. I prepared everything. Of course, the conversation didn’t go exactly the way I had expected, but I intentionally selected the words I would use in explaining things, then I became able to manipulate the target FSs naturally.

Table 5 Completing speaking journals and incorporating phrases

Item Category	N	Agree*	Mean**	Mode**	SD
5. Choosing a key word to investigate in the corpus was easy	29	9 (31%)	3.00	3	1.10
29. It was difficult to manipulate the SJ conversations to use my key words	28	27 (96%)	5.14	5	0.93
36. Context of planned conversation is important when choosing words and phrases	29	28 (97%)	4.55	5	0.91
28. It was easy to use my phrases in conversation	29	8 (28%)	2.83	2	1.07

*Raw numbers and percentage in parentheses

**1: strongly disagree, 2: disagree, 3: somewhat disagree, 4: somewhat agree, 5: agree, 6: strongly agree.

Some students were able to manipulate their FSs on the spot during their conversations, demonstrating adeptness at noticing and dominating the more open slot-and-frame patterning of some FSs. One student, searching for the word *increase* in COCA, arrived at the phrase *increase your lifespan*. However, when the time came to actually use the phrase in the real conversation, circumstances dictated that she change the phrase to *decrease your lifespan*. The student acknowledged as much in her speaking journal reflection log, demonstrating a high level of performance and agency by appropriating the phrase and manipulating it to fit her needs.

In addition to the success stories, though, Table 5 illustrates that using novel FSs in a natural way is not always easy for students. Just as learning the proper usage of novel vocabulary items can be challenging at times, it should not be too surprising that learners will occasionally encounter difficulties working prefabricated material into conversations. In their speaking journal logs, some students noted times when they abandoned target phrases because they felt unable to work them into their conversations naturally. At other times, students wrote that they were so absorbed in their conversations that they forgot to include a target phrase. Conversations are, after all, inherently open and dynamic; even the most socially alert people cannot predict with 100% accuracy how a given conversation will unfold.

Perhaps the most common difficulty throughout the course was students feeling unable to capture the nuance, or precise meaning of a target phrase, and use it in a pragmatically appropriate way in their own conversations.

One student commented in the interviews:

I believe the nuance is different when I translate the target phrases into Japanese. It seems okay to use those phrases in a straightforward way in any situation, but I was sometimes told that I couldn't use certain phrases in certain contexts because the nuance is slightly weird even though they were grammatically correct.

This comment is interesting in that corpus work lends itself to discovering more 'natural', frequently used language. While corpus-based language learning might help students discover frequently occurring sequences of words that will often sound natural in speech, we hypothesize that the situation described by the student above is the result of trying to shoehorn a more idiomatic phrase into the conversation. We conjecture that this unexpected use of idiomaticity sometimes struck the students' conversation partners as odd.

Table 6 *Appropriateness of phrases*

Raters	Appropriateness				Mean
	1	2	3	4	
A	4 (3.5%)	20 (17.5%)	36 (31.6%)	54 (47.4%)	3.23
B	18 (15.8%)	24 (21.1%)	36 (31.6%)	36 (31.6%)	2.79
C	19 (16.7%)	26 (22.8%)	18 (15.8%)	51 (44.7%)	2.89
D	6 (5.3%)	6 (5.3%)	33 (28.9%)	69 (60.5%)	3.45
Mean	11.75 (10.3%)	19 (16.7%)	30.75 (27%)	52.5 (46%)	3.09

In order to more thoroughly gauge the consistency with which students were able to nest their target phrases into a larger context in a pragmatically appropriate way, 114 items were collected and rated by four native speakers of English on a numerical rating scale measuring appropriateness, with 1 being least appropriate and 4 being most appropriate. Raters were given an Excel file with the target phrases underlined and embedded in the larger context of the conversation in one column, and a drop-down menu where they could select 1 to 4 in the adjacent column.

Consistency between raters was again calculated via Cronbach's Alpha, and was $r = .816$, indicating a high level of consistency (Table 6). The raw frequency of all phrases receiving a given rating is displayed in columns 1–4, which represent the numerical rating score, and the corresponding percentage is in parentheses. The mean score assigned by each rater is given in the last column; the overall mean score was 3.09. Also, the mean number of each score assigned by all raters is given in the bottom row. In general, the scores suggest that students were able to employ their target phrases in a pragmatically appropriate manner. However, there were still a considerable number of phrases that were used in an inappropriate manner. This is likely a reflection of the legitimacy of students' concerns over not always being able to grasp the nuance of novel FSs, and possibly a lack of sufficient planning and preparation for their real conversations.

On the other hand, it is important to note that conversations can consist of as much listening as speaking, and some students noted an increase in ability to understand their interlocutors and authentic English input, such as television programs, due to the FSs they learned through their corpus consultation. One student commented:

As I mentioned earlier, for example, the phrase, "grab a bite" is a phrase that I couldn't have encountered if I had studied in a regular way. It doesn't appear in a textbook. I find words reading news articles then I try to search for FSs around those words... When I watch TV and encounter an FS that I learned in class, I would feel "I know this meaning!"

4.5 Attitudes toward DDL

With respect to student attitudes toward DDL, there was some scepticism about the grammaticality of the concordance data, and they believed it prudent to have a dictionary on hand to verify corpus findings. This scepticism is likely why they felt that a class of this nature is better suited to advanced learners of English as opposed to novices.

Table 7 Attitudes and beliefs about data-driven learning

Item Category	N	Agree*	Mean**	Mode**	SD
33. I trust the phrases I find in the corpus to be grammatically correct	29	14 (48%)	3.48	3	1.12
32. Using a corpus is best in combination with a dictionary	29	26 (90%)	4.72	5	1.22
44. A class of this nature is better suited to advanced learners than beginners	29	25 (86%)	4.59	5	1.35
35. After taking this class I believe that grammar and vocabulary are more closely related	29	21 (72%)	4.24	4	1.21
43. I enjoyed being able to direct my own learning	29	21 (72%)	4.00	4	1.10
38. I would recommend this type of class to other English learners in Japan	27	22 (81%)	4.15	4	1.20
24. Corpus use should be taught in English classes more regularly	28	20 (71%)	3.89	4	0.99
25. I will use the corpus in future English classes	28	21 (75%)	4.11	4	0.83

*Raw numbers and percentage in parentheses

**1: strongly disagree, 2: disagree, 3: somewhat disagree, 4: somewhat agree, 5: agree, 6: strongly agree.

On the other hand, there were many positive perceptions of corpus consultation as well. Item 35 in Table 7, for example, suggests that one of the major aims of the course, to increase learner awareness of the interdependence of lexis and grammar, was largely effective with 21 of 29 participants agreeing with the statement. Additionally, the majority of students believed that they would continue to use a corpus in future classes, would recommend DDL to other learners of English in Japan, and believed that corpus use should be taught more regularly in English classes. These numbers suggest that students in this study were convinced of the utility of corpora in language education.

5 Conclusion

Chambers (2005: 111–112) notes that while there is an increasing body of research on corpus use by learners, there is “considerable scope for development, particularly in the area of course design and structure, concerning how one can successfully integrate corpus consultation into a programme of language study in higher education”. This was precisely one of the major aims of this paper. The course placed students in an interaction-rich environment that saw them interact with authentic materials in the form of news articles or videos, TV shows, etc. The students used that interaction to form a topic and choose key words, familiar or unfamiliar, that they believed would be useful in a conversation about their topic. Students then engaged in pattern hunting as they attempted to identify common patterns of usage regarding their key words. They endeavored to appropriate their newfound phrases by pushing themselves to produce them in their output in class with their classmates, and outside of class with more proficient speakers of English.

Based on the survey, interviews, and speaking journal reflection logs, students generally reported favorable impressions of the course, and perceived DDL as having a positive effect on their language learning. Participants also believed that corpus consultation should be taught more regularly in English classes, and planned to continue using the skills they

learned in the class reported on here. The usual complaints about the tedious nature of learning to use the corpus surfaced, and some participants did express reservations about being able to use their newly discovered FSs in pragmatically appropriate ways. But the sample of phrases rated in this study suggests that students were more successful than not in employing their phrases in an appropriate manner. Additionally, there is some evidence that familiarization with FSs in the course led to increased understanding of FSs encountered outside of the class. Beyond the speaking journals, students responded very positively to the student-led lessons and believe that corpora can be a good tool for discovering the difference between near-synonymous words and phrases.

There are, however, a number of limitations with this study that need to be addressed. First and foremost is the small number and homogeneity of participants. With only 29 participants, all of whom were Japanese, it is difficult to extrapolate the findings of this study to a wider array of contexts. Another serious limitation is the lack of longitudinal data. While students indicated they would continue to use the corpus-consulting skills they learned into the future and in other classes, no follow-up survey or contact was made to verify this. More longitudinal studies that track learners' corpus use over an extended period of time are a worthwhile direction for more research (cf. Yoon, 2008). It would be especially interesting to track students who go through a corpus-training course for a number of years after completion to see how long and to what extent they independently engage in corpus consultation.

This study has illustrated that, with training, learners can take advantage of the power of a corpus, and has provided qualitative evidence suggesting that students strongly believe that corpus consultation has the potential to facilitate the learning of novel usages of familiar lexical items, thus supporting the quantitative evidence provided by Frankenberg-Garcia (2012; this volume). Future research could perhaps investigate the effect of different corpus-based approaches in increasing learners' functional knowledge of familiar lexical items. For example, in addition to investigating paper-based teacher-prepared DDL activities (Boulton, 2010a; Johns, 1991), one exciting avenue could be to explore their electronic counterparts via tablet computers that offer more interactive and tactile affordances.

A number of scholars note that corpus consultation may have its brightest future outside the classroom as it affords students a high degree of autonomy (Chambers 2007; Yoon & Hirvela 2004). To see this prediction come to fruition, we recommend focusing efforts on making already excellent resources such as COCA even more accessible to casual learners. Perhaps corpora designed for hands-on use by learners can afford to sacrifice some depth and functionality in exchange for accessibility and intuitiveness. Corpus-based language learning might see even wider adoption if vetted, principled corpora were as accessible and intuitive as, say, Google searches.

References

- Aguado-Jiménez, P., Pérez-Paredes, P. and Sánchez, P. (2012) Exploring the use of multidimensional analysis of learner language to promote register awareness. *System*, **40**(1): 90–103.
- Bernardini, S. (2004) Corpora in the classroom: An overview and some reflections on future developments. In: Sinclair, J. (ed.), *How to use corpora in language teaching*. Amsterdam: John Benjamins, 15–36.
- Biber, D., Conrad, S. and Cortes, V. (2004) 'If you look at...': Lexical bundles in university teaching and textbooks. *Applied Linguistics*, **25**(3): 371–405.
- Boulton, A. (2009) Testing the limits of data-driven learning: Language proficiency and training. *ReCALL*, **21**(1): 37–54.

- Boulton, A. (2010a) Data-driven learning: Taking the computer out of the equation. *Language Learning*, **60**(3): 534–572.
- Boulton, A. (2010b) Learning outcomes from corpus consultation. In: Moreno Jaén, M., Serrano Valverde, F. and Calzada Pérez, M. (eds.), *Exploring new paths in language pedagogy: Lexis and corpus-based language teaching*. London: Equinox, 129–144.
- Chambers, A. (2005) Integrating corpus consultation in language studies. *Language Learning & Technology*, **9**(2): 111–125.
- Chambers, A. (2007) Integrating corpora in language learning and teaching. *ReCALL*, **19**(3): 249–251.
- Chambers, A. and O’Sullivan, Í. (2004) Corpus consultation and advanced learners’ writing skills in French. *ReCALL*, **16**(1): 158–172.
- Chen, Y. H. and Baker, P. (2010) Lexical bundles in L1 and L2 academic writing. *Language Learning & Technology*, **14**(2): 30–49.
- Cobb, T. and Boulton, A. (Forthcoming) Classroom applications of corpus analysis. In: Biber, D. and Reppen, R. (eds.), *Cambridge handbook of corpus linguistics*. Cambridge: Cambridge University Press.
- Conklin, K. and Schmitt, N. (2008) Formulaic sequences: Are they processed more quickly than nonformulaic language by native and nonnative speakers? *Applied Linguistics*, **29**(1): 72–89.
- Council of Europe. (2001) *Common European framework of reference for languages: Learning, teaching, assessment*. Strasbourg: Language Policy Unit. http://www.coe.int/t/dg4/linguistic/source/framework_en.pdf
- Davies, M. (2008) The Corpus of Contemporary American English: 450 million words, 1990–present. <http://corpus.byu.edu/coca/>
- Dörnyei, Z. and Taguchi, T. (2010) *Questionnaires in second language research: Construction, administration, and processing* (2nd ed.). New York: Routledge.
- Educational Testing Service. (2013) Sukoa no meyasu [Score descriptors]. <http://www.toeic.or.jp/toeic/about/result.html>
- Ellis, N. C. (2002) Frequency effects in language processing: A review with implications for theories of implicit and explicit language acquisition. *Studies in Second Language Acquisition*, **24**(2): 143–188.
- Ellis, N. C. (2012) Formulaic language and second language acquisition: Zipf and the phrasal teddy bear. *Annual Review of Applied Linguistics*, **32**(1): 17–44.
- Erman, B. and Warren, B. (2000) The idiom principle and the open choice principle. *Text*, **20**: 29–62.
- Frankenberg-Garcia, A. (2012) Learners’ use of corpus examples. *International Journal of Lexicography*, **25**(3): 273–296.
- Gries, S. T. (2008) Corpus-based methods in analyses of second language acquisition data. In: Ellis, N. C. and Robinson, P. (eds.), *Handbook of cognitive linguistics and second language acquisition*. New York/London: Routledge, 406–431.
- Gries, S. T. (2010) Behavioral profiles: A fine-grained analysis and quantitative approach in corpus-based lexical semantics. *The Mental Lexicon*, **5**(3): 323–346.
- Hoey, M. (2009) Corpus driven approaches to grammar: The search for common ground. In: Römer, U. and Schulze, R. (eds.), *Exploring the lexis-grammar interface*. Amsterdam: John Benjamins, 33–47.
- Hunston, S. (2002) *Corpora in applied linguistics*. Cambridge: Cambridge University Press.
- Irujo, S. (1986) A piece of cake: Learning and teaching idioms. *ELT Journal*, **40**(3): 236–242.
- Jiang, N. and Nekrsova, T. M. (2007) The processing of formulaic sequences by second language speakers. *The Modern Language Journal*, **91**(3): 433–445.
- Johns, T. (1986) Micro-concord: A language learner’s research tool. *System*, **14**(2): 151–162.
- Johns, T. (1991) From printout to handout: Grammar and vocabulary teaching in the context of data-driven learning. In: Johns, T. and King, P. (eds.), *Classroom concordancing*. *English Language Research Journal*, **4**: 27–45.

- Kennedy, G. (2008) Phraseology and language pedagogy: Semantic preference associated with English verbs in the British National Corpus. In: Meunier, F. and Granger, S. (eds.), *Phraseology in foreign language learning and teaching*. Amsterdam: John Benjamins, 21–41.
- Kennedy, C. and Miceli, T. (2010) Corpus-assisted creative writing: Introducing intermediate Italian learners to a corpus as a reference resource. *Language Learning & Technology*, **14**(1): 28–44.
- Liu, D. and Jiang, N. (2009) Using a corpus-based lexicogrammatical approach to grammar instruction in EFL and ESL contexts. *The Modern Language Journal*, **93**(1): 61–78.
- Mackey, A., Abbuhl, R. and Gass, S. (2012) Interactionist approach. In: Gass, S. and Mackey, A. (eds.), *The Routledge handbook of second language acquisition*. New York: Routledge, 7–23.
- Pérez-Paredes, P. and Cantos Gómez, P. (2004) Some lessons students learn: Self-discovery and corpora. In: Aston, G., Bernardini, S. and Stewart, D. (eds.), *Corpora and language learners*. Amsterdam: John Benjamins, 247–257.
- Römer, U. (2009) The inseparability of lexis and grammar. *Annual Review of Cognitive Linguistics*, **7**: 141–163.
- Römer, U. (2011) Corpus research applications in second language teaching. *Annual Review of Applied Linguistics*, **31**: 205–225.
- Sinclair, J. (1991) *Corpus, concordance, and collocation*. Oxford: Oxford University Press.
- Sinclair, J. (2004) *Trust the text: Language, corpus, and discourse*. London: Routledge.
- Swain, M. (1995) Three functions of output in second language learning. In: Cook, G. and Seidlhofer, B. (eds.), *Principles and practice in applied linguistics: Studies in honour of H. G. Widdowson*. Oxford: Oxford University Press, 125–144.
- Tremblay, A., Derwing, B., Libben, G. and Westbury, C. (2011) Processing advantages of lexical bundles: Evidence from self-paced reading and sentence recall tasks. *Language Learning*, **61**(2): 569–613.
- Wray, A. (2002) *Formulaic language and the lexicon*. Cambridge: Cambridge University Press.
- Wray, A. and Fitzpatrick, T. (2010) Pushing learners to the extreme: The artificial use of prefabricated material in conversation. *Innovation in Language Learning and Teaching*, **4**(1): 37–51.
- Yoon, H. (2008) More than a linguistic reference: The influence of corpus technology on L2 academic writing. *Language Learning & Technology*, **12**(2): 31–48.
- Yoon, H. and Hirvela, A. (2004) ESL student attitudes toward corpus use in L2 writing. *Journal of Second Language Writing*, **13**(4): 257–283.

Appendix

Interview Questions

1. For you, what are the greatest challenges in using corpora to learn English?
英語を学ぶためにコーパスを使う上で、もっとも難しかったことは何ですか。
2. What are some of the most useful and valuable things you learned in this course?
あなたがこの授業で学んだもっとも役に立った、あるいは価値のあることについて話してください。
3. Now that you have completed this course, has your view of the relationship between vocabulary and grammar changed? How?
今この授業を修了して、語彙と文法の関係性についてのあなたの考えは変わりましたか。もし変わったとすれば、どのようにですか。
4. Now that you have completed this course, has your view of the role context plays in word and phrase choice changed?
今この授業を修了して、語彙とフレーズを選ぶときのコンテキストの役割についてのあなたの考えは変わりましたか。