# Developing the preliminary essay bundles list (EBL) and its applicability to EAP

### Ryo SAWAGUCHI

# Abstract

This study developed the preliminary list of lexical bundles (e.g., on the other hand, the fact that) for argumentative essay writing and explored its potential applications to English for General Academic Purposes (EGAP) practice to prepare undergraduate students for their future use in academic written English genre (e.g., research papers). The list, called the Essay Bundles List (EBL), was created by extracting frequently used lexical bundles from opinion- and source-based argumentative essays by L1 English speakers. Corpora consulted include the International Corpus Network of Asian Learners of English (ICNALE), Louvain Corpus of Native English Essays (LOCNESS), PERSUADE2.0, Michigan Corpus of Upper-Level Student Papers (MICUSP), and British Academic Written English (BAWE). A total of 3,768 bundles were compared with the list of academic written English (Academic Formulas List: AFL) to confirm EBL applicability. The results showed that the EBL covers approximately 80% of the AFL, indicating its potential as an EGAP wordlist. Correspondence analysis of the top 21 frequent bundles in opinion- and source-based essays and the AFL revealed that the opinion-based bundles (e.g., *I believe that*) can be made suitable for academic written English with the use of inanimate subjects (e.g., it is true that), while source-based discourse bundles (e.g., in order to) imply their direct applicability. The EBL was refined to 127 bundles according to their difficulty levels on the Common European Framework of Reference (CEFR) scale (A2, B1, B2) in proficiency order. This study suggested that, basic referential bundles (e.g., the fact that) and objective stance bundles (e.g., this means that) are appropriate for A2 and B1 students. Discourse bundles (e.g., to begin with, on the one hand) should remain a focus throughout the progression from B1 to B2. Advanced referential bundles, such as the existence of, are most suitable for instruction at the B2 level.

# Introduction

Corpus linguistics developments in the last few decades have made it possible to analyze recurring word clusters called lexical bundles (e.g., in order to, as a result of). They have been examined in English for Academic Purposes (EAP) contexts (Biber et al., 2004; Biber & Barbieri, 2007; Hyland, 2008) because of their significant discourse functions in academic speech or writing. Lexical bundles are "multiword sequences that occur most commonly in a given register" (Biber & Barbieri, 2007, p. 264). Register refers to variations of language (e.g., spoken or written) used in different situational characteristics. According to Biber et al. (2004), classroom conversations are regarded as an *oral* register, while academic prose is defined as a *literate* register. Biber and Barbieri (2007, p. 273) comment that "the extent to which a speaker or a writer relies on lexical bundles is strongly influenced by their communicative purposes." For example, the bundle it is clear that expresses a writer's point of view in the following sentence, while *as a result* connects the preceding and subsequent sentences. The significance of lexical bundles has resulted in the development of a useful wordlist, the Academic Formulas List (AFL: Simpson-Vlach & Ellis, 2010), to assist in the intensive learning of lexical bundles (which researchers term *formulas*). However, no list of lexical bundles for college-level writing genres such as essays has yet been developed.

To prepare students for future academic English situations (e.g., writing papers), argumentative essays have been the most common writing genre for undergraduate students (Wu, 2006). The possible reasons could be the applicability of argumentative essays to research papers in terms of genre and text types. Swales (1990) defines genre as a set of events sharing the same communicative purposes. Biber (1989) describes text types as differences in linguistic features. In case of argumentative essays and research papers, the two express the writer's stance and support it with evidence; thus, these can be classified under the argumentative writing genre with the same communicative purpose: arguing. Johnson (2018) claims that the rhetorical characteristics in genres such as argumentation and exposition encompass various text types. Consequently, argumentative essays and research papers possibly share the similar linguistic features, e.g., the same lexical bundles. Argumentative essays have also been used to assess the use of lexical bundles by undergraduate students to improve their basic

academic writing skills (Granger, 2017; Nam & Park, 2020; Sawaguchi, 2024), a foundation for English for General Academic Purposes (EGAP; Blue, 1988) in which "students from a wide range of disciplines will write in diverse genres" (Tardy et al., 2022, p. 3). Despite the significance of argumentative essay writing in EGAP programs, no lexical bundle wordlist has yet been developed for intensive vocabulary learning for essay writing. In EAP contexts, wordlist use has facilitated greater student academic vocabulary use (Shoufaki & Petrić, 2021). Given the effectiveness of the EAP wordlist, there is a vital need for an argumentative essay writing wordlist. Another unexplored area of argumentative essay writing is its applicability to undergraduate students' future use of academic written English genre (e.g., research papers).

Therefore, developing an argumentative essay lexical bundle wordlist for academic English could significantly encourage more meaningful and focused EGAP writing practice. This study develops a possible lexical bundle list for argumentative essay writing and investigates its relevancy to academic written English. Practical suggestions on the use of the list are also proposed.

### **Literature Review**

Essays are an important text type of writing in higher education settings (Nesi et al., 2017), with argumentation in particular often being a key student requirement (Wingate, 2012) for the development of critical/logical thinking and rational argument skills.

There has been significant research into the pedagogical applications of argumentative essays. To gain a more precise understanding of argumentative essays, Yoon and Tabari (2023) classified argumentative essays into two categories: source-based and opinion-based. In source-based essays, writers organize and present their arguments based on established information sources (e.g., research articles). By contrast, opinion-based essays require the writer's knowledge or experience: the topics include the pros and cons of part-time jobs for university students.

Despite the noted importance of lexical bundles in academic English, few studies have investigated how these lexical bundle items are dealt with in argumentative essay teaching materials or curricula. Sawaguchi (2024) focuses on identifying target lexical bundles for opinion-based argumentative essay writing using the L1 English speaker

essay corpora in the Louvain Corpus of Native English Essays (LOCNESS: Granger, 1998) and the International Corpus Network of Asian Learners of English (ICNALE: Ishikawa, 2023), and identifies the most frequent lexical bundles in various opinion-based essay topics. He proposed a teaching order for the target bundles based on the estimated difficulty for Japanese university students, with a particular focus on the bundles they are unfamiliar with.

However, EAP applications must include target lexical bundles for both opinionand source-based essays, as both argumentative writing types are included in EGAP writing courses. It would also be valuable to clarify the lexical bundle differences between these two essay types. While both share the common moniker of "argumentative essays," it is likely that the lexical bundles vary because of their different argument bases: opinion or source. Understanding the lexical bundle differences for these two essay types could assist EAP teachers in teaching according to their needs. Further, because EGAP writing equips students with general academic writing skills before they proceed to English for Specific Academic Purposes (ESAP: specializing in their own disciplines; Blue, 1988), identifying the associations between argumentative essay lexical bundles and those in academic written English could also be useful for teaching practice. For example, a lexical bundle on the other hand frequently occurs in four different disciplines, namely biology, electrical engineering, applied linguistics, and business studies (Hyland, 2008). Given its cross-disciplinary usage, prioritizing on the other hand in EGAP writing courses for first- or second-year university students would help them build transferable writing skills that remain useful regardless of their specific discipline when they advance to ESAP contexts (e.g., graduate school studies) and academic writing in their respective fields. This is pertinent to "the nature of a "common core" of features relevant to all types of academic writing, applicable in a wide range of EAP teaching contexts" (Gardner et al., 2018, p. 647), which could allow students to apply their lexical bundle knowledge learned in EGAP to ESAP. The teaching practice would also be more effective if the target lexical bundles have difficulty levels suitable for students at different proficiency levels. Accordingly, this study addressed the following research questions:

RQ 1: To what extent are argumentative essay lexical bundles relevant to those in academic written English?

RQ 2: How can the argumentative writing lexical bundles be categorized

according to their difficulty levels?

The study then explored how the findings in RQ 1 and 2 could be applied to EGAP practice.

### **Data and Procedure**

### **Opinion-Based Essay Corpus**

This study used three opinion-based argumentative essay corpora: the ICNALE Written Essays (Ishikawa, 2023), the Louvain Corpus of Native English Essays (LOCNESS) (Granger, 1998), and the PERSUADE2.0 (Crossley et al., 2024). These three corpora include various topics, such as the pros and cons of animal testing and part-time jobs for university students, which are generally based on the writers' own opinions or ideas and do not typically require research-based evidence; therefore, in this paper, I termed these types of argumentative essays "opinion-based essays." The ICNALE, the LOCNESS, and the PERSUADE2.0 were chosen for the following three reasons. First, to the best of the author's knowledge, they are publicly available free corpora containing the essays by L1 English speakers, making it easier for other researchers to replicate the results of the study. Second, the corpora contain target-like lexical bundles, such as A-level essays and those written by L1 English instructors or professors. Third, the corpora have over 20 different essay topics, which allows for the extraction of commonly used lexical bundles in the corpora regardless of the topic. As in Nation's (2016) discussion on the creation of wordlists, range (See range in lexical bundle definition and extraction for details) is one of the most important criteria, as useful words should be found in a variety of texts.

During the extraction process, I excluded essays in the LOCNESS that were not argumentative, such as literary and exam essays in the file USMIXED. Table 1 presents the topics, the number of words and files for the opinion-based essay corpus, and the corpora analyzed in the study. The files part-time jobs and smoking in restaurants are from the ICNALE (the pros and cons of part-time jobs and smoking in restaurants). The topics in the LOCNESS were manually categorized into four major topics: human rights (e.g., gender equality), technology (e.g., the invention of computers), politics (e.g., parliamentary systems), and others (e.g., sports, the media). The files seeking opinions (seeking multiple opinions from others) and summer projects (should summer Table 1. Breakdown of the opinion-based essay corpus in the

~~~~		
Topics	No. of files	No. of words
Human rights	45	45,835
Politics	55	41,518
Part-time jobs	200	45,415
Seeking opinions	74	45,182
Smoking in restaurants	200	45,198
Summer projects	79	45,028
Technology	118	63,720
Other	120	88,194
Total	891	420,090

projects be designed by students?) are from the PESUADE2.0.

study

### Source-Based Essay Corpus

For the source-based essay data, I consulted the British Academic Written English (BAWE; Nesi et al., 2008), which includes course assignment essays from British university students, and the Michigan Corpus of Upper-Level Student Papers (MICUSP; Römer & O'Donnell, 2011), which contains approximately 830 A-grade papers from various disciplines (humanities and arts, social sciences, physical sciences) from the University of Michigan. Because the MICUSP and the BAWE both include research-based essays written by university students from different disciplines, these are defined as "source-based essays" in this study. As with the opinion-based essay data, only essays written by L1 English speakers were extracted. The BAWE and the MICUSP were prioritized over other similar type of corpus: Academic Writing at Ackland (AWA) due to the two corpora's potential large number of words for this study; the BAWE: approximately 580, 000 words; the MICUSP: approximately 450,000 words. These sizes of words were considered adequate for obtaining data from varied disciplines. While AWA was also a potential candidate, integrating it would have required a major adjustment to the data balance in this study. To maintain equal representation across disciplines, the study focused on gathering an equal number of words from the arts and humanities (including social sciences) and the sciences. Given the aim of developing an EGAP essay lexical bundle list applicable across disciplines, the dataset was structured to ensure balanced discipline coverage. Table 2 shows the

discipline genres and the total number of files and words analyzed in this study. The arts and humanities were subdivided into specific disciplines, such as archeology, linguistics, and history, and the social sciences had disciplines including business, economics, and education. Compared to the arts and humanities and social sciences, both the BAWE and the MICUSP have a relatively limited number of essays from the life/physical sciences (e.g., biology, physics). Therefore, life and physical sciences were integrated into the sciences to balance the number of words reviewed in both the arts and humanities and sciences to approximately 340,000 words each.

Table 2. Breakdown of the source-based essay corpus in the study

Disciplines	No. of files	No. of words
Arts and humanities	117	348,379
Social sciences	136	345,859
Sciences (life/physical)	171	344,343
Total	424	1,038,581

# Lexical Bundle Definition and Extraction

This study defines lexical bundles as three-to-five-word clusters that satisfy the following frequency and range criteria. The reason I focused on three to five clusters is discussed first.

Word cluster length: The RQ 1 of this study is focused on exploring how the lexical bundles in the essay wordlist could be applied to academic written English. To do this, I examined the coverage of the essay wordlist in the AFL, for which I decided the lexical bundle lengths should be the same. For example, the bundle *at the end of the day* is a six-word bundle found in the essay wordlist; however, the AFL limits bundle lengths to five, which results in bundles such as *the end of the*. By setting an equal length for the word clusters, this study sought to discover the bundles that overlap the argumentative essay bundles and the AFL.

Frequency: Lexical bundles occur at least 20–40 times per million or more in different texts (Biber & Barbieri, 2007; Hyland, 2008). These frequency criteria indicate that lexical bundles do not occur by chance but are a representation of linguistic phenomena. This study employed the standard 20 times per million criteria for the extraction of both the opinion- and source-based lexical bundles. Compared to

studies such as Hyland (2008), which analyzed 3 million words, this study used relatively small-sized corpora (approximately 1.45 million words in total), primarily because the study sought to identify the frequently appearing lexical bundles in smaller L1 corpora by establishing a minimum frequency threshold.

Range: Range is the extent to which lexical bundles are distributed across various texts. Research has indicated that lexical bundles appear in five or more texts (subcorpora created from the main corpus) (Biber et al., 2004; Bychkovska & Lee, 2017; Omidian et al., 2018). Range is a key criterion for filtering an individual writer's idiosyncratic language use. For example, if one writer uses the bundle as a result three times, this would affect the total number of raw frequencies; however, analyzing texts by different writers reduces this risk. I applied different range criteria for the opinionbased and source-based lexical bundle extractions. For the opinion-based bundles, I set a minimum of three different texts because the size of the opinion-based essay corpora in this study was similar to that consulted in Chen and Baker (2016), who set three ranges and analyzed under 1 million words (approximately 200,000 words). By setting a lenient range criterion, this study gathered as many lexical bundles as possible from the relatively small-sized opinion-based essay corpora. For the source-based lexical bundles, I applied five different text criteria because this was the standard criteria in previous studies; Omidian et al (2018), whose corpus size was very similar to this study (1030,000 words), used a five-range criterion.

All extraction processes were performed using the N-gram function in the computer concordance software AntConc Ver. 4.2.4 (Anthony, 2023). The extraction resulted in 3,768 opinion- and source-based lexical bundles. However, among the bundles that met the aforementioned frequency, range, word length criteria but are strongly topic-related bundles such as *part-time jobs* in the topic part-time jobs for university students were manually excluded from the analysis because of their low pedagogical value for essay writing. Specifically, *part-time jobs* had the highest frequency (943 times per million words) followed by *be able to* (567 times per million words) in the three-word bundles in opinion-based corpus. Despite the high frequency of *part-time jobs*, the bundle was excluded from the analysis involving frequency information. In contrast, the bundles consisting solely of function words (e.g., *this is a*) were included for the analysis in accordance with the criteria employed in the AFL, which regards these as lexical bundles. Hereafter, the bundles list is called the Essay

### Bundles List (EBL).

### **Compatible Academic Written£ English**

For comparison purposes, this study termed the AFL (Simpson-Vlach & Ellis, 2010) as "academic written English". The AFL is the largest wordlist to date that contains academic English lexical bundles (e.g., *on the other hand, as a result of*) commonly used across disciplines (e.g., social sciences, humanities, medicine) whose coverage facilitated comparison with this study's aim to develop an EGAP wordlist that could be applicable regardless of disciplines. Another advantage of the AFL is that it categorizes the lexical bundles into three major functional categories: referential (e.g., *in the case of*), stance (e.g., *it is important to*), and discourse (e.g., *in order to*), which allowed for in-depth interpretations of the similarities and differences between the essay lexical bundles in the study and those in the AFL in terms of discourse functions.

The AFL has both spoken and written academic lexical bundle lists termed as written/spoken AFL respectively. Written AFL consists of the lexical bundles frequent in academic written English text types (e.g., research papers, textbooks), while spoken AFL includes the frequent bundles in spoken academic English registers (e.g., lectures, seminars). The AFL integrates these bundles to the *core* AFL, whose lexical bundles are commonly used in both academic speech and writing. Since opinion-based lexical bundles are often more colloquial (Chen & Baker, 2016), this study chose the core AFL to better assess its coverage in the EBL. Additionally, the core AFL contains more lexical bundles (207) compared to the spoken and written AFL (200 bundles each), making it more extensive for the coverage assessment of the study, which involves a total of 3,768 lexical bundles in the EBL. While the core AFL provides frequency information for both spoken and written academic English, this study focused on the frequency data for written academic English to ensure a consistent comparison with the EBL. Hereafter, the core AFL will be simply termed as AFL.

# **Results and Discussion**

#### The Applicability of the EBL to Academic Written English

RQ 1 of the study explored the applicability of the EBL to academic written English. For this purpose, the coverage (matching rate) of the EBL and the core AFL lexical bundles was investigated. Table 3 shows the EBL coverage in the AFL and

reveals that all lexical bundles in the EBL overlapped 78.7 % of the total 207 lexical bundles in the AFL, which indicates that the EBL has a high degree of coverage in academic written English, and significant potential for inclusion in an EGAP wordlist to prepare students for the future use of academic written English.

Table 3. Coverage of the EBL in the AFL

EBL overlapping bundles	AFL bundles	Coverage (%)
163	207	78.7%

To further explore the frequency relationship between the argumentative (opinion- and source-based) bundles and the AFL (written academic English), a correspondence analysis was conducted on the top 21 frequent AFL (top 10% of the 207 AFL) and the EBL (source/opinion) corresponding 21 bundles using the langtest.jp (Mizumoto, 2015), which is a multifunctional application website that performs statistical analyses. Figure 1 shows the biplot of the correspondence analysis. Dimen-



Figure 1. Correspondence analysis of the opinion, source, and AFL bundles

sion 1 (horizontal line) and 2 (vertical line) have the following eigen values and contribution rates: Dimension 1: eigen value 0.27, contribution rate 88.7%; Dimension 2: eigen value 0.03, contribution rate 11.3%. The column scores (locations on the biplot) for opinion, source, and AFL are as follows:

Opinion: Dimension 1 = -1.08, Dimension 2 = 0.33

Source: Dimension 1 = 0.50, Dimension 2 = -1.37

AFL: Dimension 1 = 1.30, Dimension 2 = 1.21.

Table 4 presents the noticeable bundles and their row scores of opinion, source, and AFL.

Bundles	Dimension 1	Dimension 2
I believe that	-2.08	1.81
the presence of	1.54	1.81
the importance of	1.13	-1.69
in order to	0.27	-1.09
the fact that	0.14	-0.33

Table 4. Characteristic bundles and row scores of opinion, source, and AFL

Figure 1 shows that dimension 2 (vertical line) separates the opinion-based bundles from the source-based and AFL bundles. One feature of opinion-based bundles is that they are characterized by assertive stance bundles e.g., *I believe that*, as shown in the upper left (dimension 1: -2.08; dimension 2: 1.81) in Figure 1. Because *I believe that* is never used in the AFL, some opinion-based stance bundles are too subjective for academic written English. Meanwhile, the lower right in Figure 1 demonstrates that the stance bundle *the importance of* (dimension 1: 1.13; dimension 2: -1.69) is frequent in source-based essays. This highlights an interesting difference between stance bundles in source-based and opinion-based essays; source-based essays take an objective stance with an inanimate subject *the importance*, while opinion-based essays display a subjective stance with a personal subject *I*. This could be due to the source differences the two argumentative essays base their arguments on; opinion: the writer's opinion or knowledge, source: objective evidence such as research articles. The similarity in frequency between *the importance of* with AFL (located on the right of dimension 1) suggests that the academic tone of source-based essays is closer to AFL. This aligns

with Granger (2017), who found that noun-based bundles are a feature of academic writing. Another similarity of the source-based bundles with the AFL is the frequency of discourse bundles such as *in order to*. While this discourse bundle is located slightly on the lower right (dimension 2: -1.09), which shows the bundle's specificity to source-based essays, it has the potential applicability to academic written English, as it is also placed on the right of dimension 1. A moderate correlation (r = .60) of the top 21 source-based bundles with those in the AFL in frequency also reinforces their potential utility.

The upper right in Figure 1 implies that the AFL is distinguished by more objective noun-based bundles (e.g. *the number/presence of*) than opinion/source-based essays that have noticeable stance bundles such as *I believe that* and *the importance of*. This difference in argumentative tone should be considered in the applications of essays to academic writing.

Placed near the center in Figure 1 (dimension 1:0.14; dimension 2: -0.33), the referential bundle *the fact that* is commonly used regardless of text types (essays and research papers). This implies that *the fact that* is an objective and widely applicable academic bundle, which makes it an essential focus in the early stages of EGAP writing instruction.

In sum, the correspondence analysis revealed that (1) opinion-based bundles, especially stance (e.g., *I believe that*) ones, are too subjective and may not be suitable for academic written English ; (2) source-based bundles are more similar to academic written English than opinion-based bundles, as shown in the high frequency of discourse bundles such as *in order to*, *as well as*, and (3) referential (e.g., *the fact that*) in argumentative essays are widely applicable to academic written English.

To gain deeper insights into the bundle match rates and detailed frequency information, Table 5 provides the top 21 frequent bundles of the EBL (opinion/source) and the corresponding AFL bundles by frequency per million words.

The frequency information of each bundle in Table 5 strengthens the points discussed in the result of the correspondence analysis. First, the top two discourse bundles (*in order to, as well as*) in source-based essays bear a strong similarity with those in the AFL. Interestingly, the two bundles exhibit the same frequency order in the source and the AFL, with *in order to* being followed by *as well as*. The prominence of these two discourse bundles illustrates one feature of academic writing, which utilizes

Opinion	Freq.	Source	Freq.	AFL	Freq.
be able to	567	in order to	487	in terms of	282
I think that	524	as well as	411	the use of	270
that it is	495	the fact that	298	in order to	255
a lot of	483	one of the	295	as well as	255
one of the	412	the use of	262	the number of	246
to have a	390	in terms of	260	there is a	223
in order to	388	there is a	233	part of the	216
the fact that	329	due to the	222	a number of	215
it would be	283	that it is	216	the fact that	203
it is a	276	as a result	207	it is not	188
as well as	269	on the other	180	there is no	185
it is not	252	such as the	178	the case of	168
there is no	248	it is not	173	in which the	166
there is a	243	part of the	168	in the case	153
I believe that	243	be able to	167	in the case of	135
the right to	243	the other hand	160	based on the	134
should not be	233	on the other hand	159	the presence of	130
that they are	217	the importance of	156	due to the	127
this is a	212	a number of	154	as a result	125
because of the	202	the development of	154	the development of	121
in the world	198	there is no	152	the role of	121

Table 5. Top 21 lexical bundles in opinion, source, and the AFL

Italic = shared in opinion and source, shading = shared in source and AFL, **bold** = shared in all the three

the two bundles to create or organize logical connections of information. *In order to* also frequently appears in opinion-based essays, which suggests the bundle's adaptability to academic writing. Another similarity between source-based bundles and those in the AFL is the frequency of noun-based bundles. *In terms of* and *the use of* are ranked within the top 6 in both source-based essays and the AFL. This again demonstrates the high applicability of source-based bundles. Second, the referential bundle *the fact that* is shared in all the three (opinion, source, and AFL), meaning it is a common bundle applicable to a range of academic writing texts.

## Dividing the EBL According to Difficulty Level

RQ 2 of the study examined the possible divisions of the EBL (3, 768 words) to

facilitate its use in EGAP writing practice. As Nation (2016) pointed out, wordlists with numerous numbers of words (e.g., 1,000 words long) are too extensive to incorporate into a particular curriculum or course. Consequently, this study classified the EBL based on difficulty level of each bundle.

### The CEFR and English Vocabulary Profile

To gather information on the bundles' difficulty level, this study referred to the Common European Framework of Reference for Languages (CEFR; Council of Europe, 2001) and the English Vocabulary Profile (EVP; Capel, 2015). The CEFR categorizes foreign language learners' proficiency into six levels: beginner (A1), elementary (A2), intermediate (B1), upper-intermediate (B2), advanced (C1), and proficiency (C2), with A1 being the lowest and the C2 being the highest. The EVP utilizes actual learner-produced data (essays) to offer CEFR-based difficulty levels for phrases (lexical bundles). For example, the EVP states that A2 learners are expected to productively use the bundle *it is true that* in writing; thus, the bundle is at A2 level.

The three-to-five 3,768 bundles in the EBL were manually checked with the corresponding CEFR levels in the EVP. For appropriate difficulty levels, the classification was limited to A2 (elementary), B1 (intermediate), and B2 (upper-intermediate). This aligns with previous studies that focused on the Asian university students at these levels, including Japanese (Nam & Park, 2020; Sawaguchi, 2024). Table 6 shows the EBL divided into A2, B1, and B2 levels.

CEFR level	No. of bundles	Proportions (%)
A2	20	15.7%
B1	40	31.4%
B2	67	52.7%
Total	127	99.8%

Table 6. CEFR-labelled EBL bundles

Note: Percentages may not sum to exactly 100% due to rounding.

As shown in Table 6, B2 (upper-intermediate) level bundles occupy the largest proportion of the labelled CEFR levels. In the previous studies that targeted Asian university students (Nam & Park, 2020; Sawaguchi, 2024), B2 level students are considered the most proficient. This suggests that, overall, the EBL has challenging learning items for average Japanese university students. These include the bundles with

relatively advanced vocabulary, including *the distinction between* whose content word *distinction* is at 5,000 level in the New Word Level Checker (NWLC; Mizumoto, 2021). The abundance of B2 bundles in source-based essays contributes to the overall large number of B2 bundles (52 of 67). In comparison, A2 (elementary) bundles account for the smallest proportion of the total labelled bundles (20 of 127). Some of them are characterized by basic vocabulary (e.g., *the fact that, it is true that*).

The proportional features of B1 and B2, which account for over 80% of the total A2, B1, and B2 bundles show the overall tendency that the EBL frequently employ the bundles that help university students present clear and logical arguments on various topics. This competency is in line with the CEFR descriptors of B1 "give reasons and explanations for opinions" and B2 "produce clear, detailed text on a wide range of subjects" (Council of Europe, 2001, p. 24), which again enhances the EBL's value to improve university students' basic academic writing skills.

### Applying the EBL to the EGAP Writing Practice

Building on the CEFR categorization, the result of RQ 1 (discourse functions and the frequency of the EBL bundles and their similarities to the AFL), and the relevant findings in previous studies, I will discuss the applications of the EBL to EGAP essay writing activities. Model answer sentences were generated by ChatGPT 4o, and later modified by the author. Figure 2 illustrates how opinion-based bundles can be used in

#### The Debate on Free University Education: A Policy Worth Considering

Introduction: Whether or not free university education is a viable policy remains a contentious issue.... it is true that the financial cost of such a policy would be enormous. However, this **does not mean**... Body 1: To begin with, free university education would lead to a more equitable society....the fact that some European countries, such as Germany and Sweden, have already implemented free university education....

Body 2: The amount of money needed to sustain free university education is another important consideration...

Conclusion: All in all, it is clear that such a policy would have positive effects....

Answer: A2: *it is true that, the fact that* B1: *whether or not, does not mean* B2: *to begin with, all in all* 

teaching A2 and B1 university students. The answers with CEFR levels are also provided.

As discussed in RQ 1, *the fact that* is a commonly used academic referential bundle. It would be effective to focus first on the bundle. It is at A2 level in the EVP, whose literal meaning and lower level of vocabulary *fact* would facilitate A2 students' use of the bundles. The stance bundles *it is true that* and *does not mean* would also be useful to develop strong arguments in writing. As discovered in RQ 1, stance bundles with personal subjects (e.g., *I believe that*) are too subjective for academic writing; thus, using inanimate subjects *it* and *this* as in Figure 2 assists in maintaining objective academic tone. In fact, the *it is* construction is frequently used in the AFL (e.g., *it is important/necessary/possible* to...). Regarding discourse bundles (*whether or not, to begin with, all in all*), all of them are ranked at B1 or B2 in the EVP. These bundles can be considered appropriate difficulty for B1 and B2 learners. These discourse bundles in opinion-based essays (e.g., *in order to*) show a high frequency similar to that of AFL (ranked within the top 21).

The above suggestions for bundles in terms of difficulty and discourse functions are also supported by previous studies. Chen & Baker (2016) found that B2 students use more objective stance bundles with *it is* constructions, and Sawaguchi (2024) discovered the B1 students' competency development to employ varied discourse bundles (e.g., *it is up to*) compared to A2 students. Opinion-based stance bundles with inanimate subjects (e.g., *it is true that, does not mean*) are beneficial for A2 students to be aware of academic stance tone at the early stages of writing practice. B1 students can also increase their repertoire of discourse bundles with the focus on those at B1 and B2 levels (e.g., *whether or not, to begin with, all in all*).

Figure 3 presents the application of source-based bundles for B2 students.

It was found in RQ 1 that source-based essays frequently employ discourse bundles (e.g., *in order to*), which facilitate the organization of presenting information. For B2 students, continued focus on formal B2 discourse bundles like the ones in Figure 3 (*despite the fact that, one the one hand*) will assist B2 students in presenting their arguments more logically, because the two bundles contrast both sides of arguments in an objective manner. At B2 level, the effective use of advanced vocabulary is also necessary. As Figure 3 shows, the referential bundles with relatively advanced

#### Renewable Energy versus Fossil Fuels: A Choice for the Future

Introduction: As climate change impacts grow, the debate over energy sources is more urgent than ever. **Despite the fact that** fossil fuels have been central to global energy, renewable sources are increasingly viewed as essential for sustainability.

Body 1: **On the one hand**, fossil fuels have been deeply embedded in global economies for centuries, but on the other,...

Body 2: However, the existence of renewable energy in our current energy mix is still limited.

...While fossil fuels and renewable energy are often compared on environmental grounds, **the distinction between** the two also lies in their economic implications....**The origins of** fossil fuels...

Conclusion: In conclusion, while fossil fuels remain essential to the global energy supply, renewable energy can be a viable and sustainable energy source for the future.

Answer: B2: *despite the fact that, one the one hand, the existence of, the distinction between, the origins of Note:* References are required in actual source-based argumentative essays.

Figure 3. Application of source-based bundles

levels of content words (*the existence of, the distinction between, the origins of*) are at B2 level in the EVP. One feature of academic written English (AFL) is the frequent use of various referential bundles, including *the presence/development of*. Aiming at the referential bundles such as *the existence of, the distinction between*, and *the origins of* will further increase B2 students' use of sophisticated referential bundles.

# Conclusion

This study sought to develop an initial framework for the list of lexical bundles for argumentative essay writing and to explore the list's potential applications to EGAP practice.

RQ 1 found that approximately 80 % of the lexical bundles in the EBL overlapped with those in the AFL, which suggests the potential for the application of argumentative essay writing to EAP. The analyses of the highly frequent top 21 bundles in the EBL and the AFL revealed the following: Opinion-based essays contain remarkable stance bundles such as *I believe that*, which may not be used in academic written English practice due to their subjectivity, while the objective referential bundles including *the fact that* is applicable. In contrast, source-based essays have the abundant discourse bundles (e.g., *in order to, due to the*), which are more similar to academic written English. Academic written English (AFL) is distinguished from both types of essays, with more frequent use of noun-based referential bundles (e.g., *the number/ presence of*).

RQ 2 classified the EBL into the CEFR-based difficulty (A2, B1, B2) level, suggesting appropriate bundles to teach at each proficiency level. Specifically, the basic referential bundles (e.g., *the fact that*) and objective stance bundles (e.g., *it is true that*) can be appropriate at A2 and B1 levels; discourse bundles (e.g., *to begin with, on the one hand*) should be the continued focus from B1 to B2 levels. Advanced referential bundles such as *the existence of* can be taught at B2 level.

Finally, the limitations of this study and the directions for future research are discussed. While the findings highlight the relevance of the EBL to EGAP instruction, further validation and adjustments are needed to refine the list and confirm its pedagogical effectiveness; thus, the list will not be publicly released at this stage. As this study represents the first attempt to develop a collection of essay-specific bundles, the study serves as a foundation for future research on the practicality of the EBL in EGAP, contributing to the development of argumentative and academic writing instruction.

### Acknowledgements

This study was supported by JST SPRING (grant number JPMJSP2150).

### References

- Anthony, L. (2023). AntConc (Version 4.2.4) [Computer Software]. Tokyo, Japan: Waseda University. https://www.laurenceanthony.net/software
- Biber, D. (1989). A typology of English texts. *Linguistics*, 27(1), 3–43.
- Biber, D., Conrad, S., & Cortes, V. (2004). "If you look at ...": Lexical bundles in university teaching and textbooks. *Applied Linguistics*, 25, 371–405.
- Biber, D., & Barbieri, F. (2007). Lexical bundles in university spoken and written registers. English for Specific Purposes, 26, 263–286.
- Blue, G. (1988). Individualising academic writing tuition. In P. C. Robinson (Ed.), Academic writing: Process and product (ELT Documents 129, pp. 95–99). Modern English Publications.
- Bychkovska, T., & Lee, J. J. (2017). At the same time: Lexical bundles in L1 and L2 university

student argumentative writing. Journal of English for Academic Purposes, 30, 38-52.

- Capel, A. (2015). The English vocabulary profile. In Harrison, J., & Barker, F. (Eds), *English profile in practice*. Cambridge University Press.
- Chen, Y. H., & Baker, P. (2016). Investigating criterial discourse features across second language development: Lexical bundles in rated learner essays, CEFR B1, B2 and C1. *Applied Linguistics*, 37(6), 849–880.
- Cortes, V. (2004). Lexical bundles in published and student disciplinary writing: Examples from history and biology. *English for Specific Purposes*, 23, 397–423.
- Council of Europe. (2001). *The common European framework of reference for languages: Learning, teaching, assessment.* Cambridge University Press.
- Crossley, S.A., Tian, Y., Boffour, P., Franklin, A., Benner, M., & Boser, U. (2024). A largescale corpus for assessing written argumentation: PERSUADE 2.0. Assessing Writing, 61, Article 100865.
- Gardner, S., Nesi, H., & Biber, D. (2018). Discipline, level, genre: Integrating situational perspectives in a new MD analysis of university student writing. *Applied Linguistics*, 40(4), 646–674.
- Granger, S. (1998). The computer learner corpus: A versatile new source of data for SLA research. In Granger, S. (Ed.), *Learner English on computer* (pp. 3–18). Addison Wesley Longman.
- Granger, S. (2017). Academic phraseology: A key ingredient in successful L2 academic literacy. In Vatvedt Fjeld, R., Hagen, K., Henriksen, B., Johannson, S. Olsen, S., & Prentice, J. (Eds.), Academic Language in a Nordic Setting: Linguistic and Educational Perspectives. Oslo Studies in Language, 9(3), 9–27.
- Hyland, K. (2008). As can be seen: Lexical bundles and disciplinary variation. *English for Specific Purposes*, *27*(1), 4–21.
- Ishikawa, S. (2023). The ICNALE guide: An introduction to a learner corpus study on Asian learners' L2 English. Routledge.
- Johnson, D. (2018). Teaching English for academic purposes in New Zealand: Making sense of genre-based instruction. In Wong, L.T., & Wong, Heidi. W. L (Eds.), *Teaching and learning English for academic purposes: Current research and practices* (pp. 239–253). Nova Science Publishers.
- Mizumoto, A. (2015). Langtest. (Version1.0) [Web application]. https://langtest.jp/
- Mizumoto, A. (2021). New Word Level Checker [Web application]. https://nwlc.pythonanywhere.com.
- Nam, D., & Park, K. (2020). Lexical bundles as criterial features in L2 academic writing: structural differences between CEFR A2 and B2 essays. *Multimedia-Assisted Language Learning*, 23(3), 68–86.
- Nation, I.S.P. (2016). *Making and using wordlists for language learning and testing*. John Benjamins Publishing company.

- Nesi, H., Gardner, S., Thompson, P., & Wickens, P. (2008). British academic written English corpus (BAWE). Centre for Corpus Research, University of Birmingham.
- Nesi, H., & Matheson, N., & Basturkmen, H. (2017). University literature essays in the UK, New Zealand and the USA: Implications for EAP. New Zealand Studies in Applied Linguistics, 23(2), 25–38.
- Omidian, T., Shafriari, H., & Siyanova-Chanturia, A. (2018). A cross-disciplinary investigation of multi-word expressions in the moves of research article abstracts. *Journal of English for Academic Purposes*, 36, 1–14.
- Römer, U., & O'Donnell, M. B. (2011). Michigan corpus of upper-level student papers (MICUSP). University of Michigan, Ann Arbor.
- Sawaguchi, R. (2024). Potential of L1 and L2 corpora to identify target lexical bundles for argumentative essay writing. *Asia Pacific Journal of Corpus Research*, 5(1), 1–21.
- Shoufaki, S., & Petrić, B. (2021). Academic vocabulary in an EAP course: Opportunities for incidental learning from printed teaching materials developed in-house. *English for Specific Purposes*, 63, 71–85.
- Simpson-Vlach, R., & Ellis, N. C. (2010). An academic formulas list: New methods in phraseology research. *Applied Linguistics*, *31*(12), 487–512.
- Swales, J.M. (1990). *Genre analysis: English in academic and research settings*. Cambridge University Press.
- Tardy, C. M., Buck, R., Jacobson, B., LaMance, R., Pawlowski, M., Slinkard, J. R., & Vogel, S.
  M. (2022). "It's complicated and nuanced": Teaching genre awareness in English for general academic purposes. *Journal of English for Academic Purposes*, 57, 101117.
- Wingate, U. (2012). 'Argument!' helping students understand what essay writing is about. Journal of English for Academic Purposes, 11(2), 145–154.
- Wu, S. M. (2006). Creating a contrastive rhetorical stance: Investigating the strategy of problematization in students' argumentation. *REC journal*, 37(3), 329–353.
- Yoon, H.J., & Tabari, M.A. (2023). Authorial voice in source-based and opinion-based argumentative writing: Patterns of voice across task types and proficiency levels. *Journal* of English for Academic Purposes, 62, Article 101228.

(Ryo Sawaguchi, Graduate School of Foreign Language Education and Research, Kansai University: rswgch@gmail.com)