



**UNIVERSITY OF LEEDS**

This is a repository copy of *The Academic Spoken Word List*.

White Rose Research Online URL for this paper:  
<http://eprints.whiterose.ac.uk/135479/>

Version: Accepted Version

---

**Article:**

Dang, TNY [orcid.org/0000-0002-3189-7776](https://orcid.org/0000-0002-3189-7776), Coxhead, A and Webb, S (2017) The Academic Spoken Word List. *Language Learning*, 67 (4). pp. 959-997. ISSN 0023-8333

<https://doi.org/10.1111/lang.12253>

---

© 2017, Language Learning Research Club, University of Michigan. This is an author produced version of a paper published in *Language Learning*. Uploaded in accordance with the publisher's self-archiving policy.

**Reuse**

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.



[eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk)  
<https://eprints.whiterose.ac.uk/>

# THE ACADEMIC SPOKEN WORD LIST

**Thi Ngoc Yen Dang**

**Averil Coxhead**

**Stuart Webb**

**Dang, T. N. Y.,** Coxhead, A., & Webb, S. (2017). The academic spoken word list. *Language Learning*, 67(4), 959–997.

The linguistic features of academic spoken and written English are different (Biber, 2006). An Academic Spoken Word List (ASWL) was developed and validated to help second language (L2) learners enhance their comprehension of academic speech in English-medium universities. It contains 1,741 word-families with high frequency and wide range in a 13-million running-word academic spoken corpus. The ASWL represents the vocabulary from 24 subjects across four equally-sized disciplinary sub-corpora. Its coverage in academic speech, academic writing, and non-academic speech indicates that the ASWL truly represents the most-frequent and wide-ranging words in academic speech. The list is graded into four levels according to Nation's (2012) BNC/COCA lists. Each level is divided into sub-lists of function words and lexical words. Users can choose the words from the list that are most suitable for learning. Depending on their vocabulary levels, learners may reach 92%-96% coverage of academic speech with the aid of the ASWL.

**Key words:** English for academic purposes; vocabulary; corpora; English as a Foreign Language; TESOL; academic spoken discourse

## Introduction

To achieve academic success at English-medium universities, second language (L2) learners need to comprehend reading materials, lectures, seminars, labs, and tutorials (Becker, 2016; Biber, 2006). These speech events are essential components of university study (Lynch, 2011). Yet, comprehending academic spoken English is challenging for L2 learners in different contexts (Flowerdew & Miller, 1992; Mulligan & Kirkpatrick, 2000). Insufficient vocabulary knowledge is frequently cited as a major reason for this difficulty (Berman & Cheng, 2001; Flowerdew & Miller, 1992). Because vocabulary knowledge and listening comprehension are closely related (van Zeeland & Schmitt, 2013), it is crucial for L2 learners to master the words that they are

likely to encounter often in a wide range of academic speech. Unfortunately, there is a lack of research in this area. Because English is a medium of instruction at tertiary levels, face-to-face and through distance learning and online open courses in English-speaking and non-English speaking countries, the demand to master academic vocabulary comes from L2 learners in a wide range of contexts with different proficiencies. To meet this demand, the determination of the most important words in academic spoken English for L2 learners should take learners' proficiencies into account.

The present study had three aims. The first aim was to identify which items appeared with high frequency in spoken discourse in a wide range of academic subjects and include them in an Academic Spoken Word List (ASWL). The second aim was to see if the ASWL truly reflects academic spoken language. The final aim was to determine the potential coverage that learners with different proficiencies may reach if they learn the ASWL. Overall, the research sought to provide a general academic spoken wordlist that is useful for L2 learners in English for General Academic Purposes (EGAP) programs regardless of their subject areas and proficiency levels.

### **Is there a need for general academic wordlists?**

An important issue when developing wordlists to help L2 learners improve their comprehension of academic spoken English is the question of whether there exists a core vocabulary in academic English. There are two different views towards this issue. One view suggests that there is a core vocabulary across multiple academic disciplines, and supports the development of general academic wordlists for L2 learners irrespective of their academic disciplines (e.g., Coxhead, 2000; Gardner & Davies, 2014; Xue & Nation, 1984). The second view questions the existence of a core academic vocabulary among different academic disciplines and argues that frequency, range, meanings, functions, and collocations of a certain word change across disciplines due to variations in the practice and discourse of disciplines (Hyland & Tse, 2007). This view promotes the idea of developing discipline-specific wordlists.

More recently, Hyland (2016) points out that the general and specific EAP approaches should be seen as ends of a continuum rather than a dichotomy. This means specificity in wordlist construction should be implemented with flexibility and consideration of the circumstances of particular students in a class. Depending on the particular teaching and learning context, one kind of wordlist may be more suitable than the other.

In English for Specific Academic Purposes (ESAP) or English for Specific Purposes (ESP) programs, where learners have highly specific needs, and plan to study similar subject areas (e.g., Mathematics, Physics, Chemistry) or even the same subject area (e.g., Mathematics), discipline-specific wordlists may better serve learners' needs than general academic wordlists. Specialized vocabulary tends to occur more often in specialized texts (Chung & Nation, 2004; Nation, 2016). Compared with general academic wordlists, discipline-specific wordlists, are better at drawing their attention to the most frequent and wide ranging words in their specific areas and providing a shortcut to reduce the amount of learning (Nation, 2013). Moreover, learners might be motivated to learn items from these lists because they can clearly see the relationship between what they study in their English courses and their subject courses (Hyland, 2016). Additionally, similarities in the learners' academic discipline may make it easier for teachers to focus on more specialized vocabulary in that discipline.

General academic wordlists have a wider application (Hyland, 2016; Nation, 2013). They can be useful in EGAP programs where learners (1) are more heterogeneous in terms of disciplines that they plan to study, (2) have not yet identified their target disciplines, (3) plan to study interdisciplinary subject areas, or (4) teachers lack background knowledge of learners' specific disciplines. In such environments, it is usually challenging for EAP teachers to satisfy the specific needs of every learner in their programs and a general academic wordlist would therefore be more practical. The value of general academic wordlists is evident from Banister's (2016) finding that Coxhead's (2000) Academic Word List (AWL), a general academic wordlist, was widely used and perceived as a useful instrument for L2 learners from a wide range of subjects by EAP teachers.

Moreover, students who begin university studies do not only take courses within a single subject area (Coxhead & Hirsh, 2007). First and second year courses are often from multiple subject areas (e.g., Engineering, Mathematics, Physics). In this respect, focusing L2 learners' attention on the shared items between multiple academic subjects may have great value at least in the initial years of study. It allows learners to comprehend the words in a range of disciplines and contexts (Nation, 2013).

One common criticism of general academic wordlists is that drawing learners' attention to the core vocabulary may neglect the discipline-specific meanings of the words. Nation (2013),

however, points out that the core meaning and discipline-specific meanings should not be seen as different from each other. Knowledge of the core meaning provides an excellent scaffolding for the acquisition of discipline-specific meanings (Crossley, Salsbury, & McNamara, 2010). Highly frequent meanings are more likely to be stored as separate entries in the brain while less frequent meanings are more likely to be inferred from the context. Therefore, knowledge of the core meaning of an academic word will help learners to gradually become aware of its discipline-specific meanings if they meet the word very often in texts from their specific subject areas. These multiple encounters of the words from general academic wordlists in different contexts related to their target subject areas help to enrich learners' knowledge of the specific senses and store them in their brains.

### **What is the value of corpora in developing general academic wordlists for L2 learners?**

Although the relative value of a word for L2 learners depends on many factors, frequency of a word in actual language use is an important one. According to Zipf's (1935) law, in a collection of texts representing a certain discourse type, the majority of words occur very infrequently; yet, a small number of words occur very frequently. Helping learners master words in the latter group is beneficial because they only need to study a small number of items in order to recognize a larger proportion of words in that kind of discourse (Nation, 2016).

A popular way to develop wordlists for L2 learners is to use corpora. Data from large and representative corpora capture actual language use, and provide a powerful and reliable way to identify the most frequent and wide ranging words in academic discourse for EGAP learners (O'Keeffe, McCarthy, & Carter, 2007). For example, Coxhead's (2000) AWL was developed from a 3.5-million word academic written corpus. Learning the 570 AWL words may allow learners to recognize around 10% of the words in a wide range of academic writing (e.g., Coxhead, 2011). Intervention studies (e.g., Townsend & Collins, 2009) showed that the AWL helps L2 learners improve their comprehension of academic written texts and academic achievement. The AWL has been widely used in EAP teaching materials design, vocabulary tests, and dictionaries (Coxhead, 2011, 2016). Another example is Gardner and Davies's (2014) Academic Vocabulary List (AVL), which was derived from a 120-million word academic written corpus. Studying the 3,000 AVL lemmas may enable learners to recognize approximately

14% of the words in academic written English. The AVL website (<http://corpus.byu.edu/coca/>) is a valuable vocabulary resource for researchers, teachers, and learners.

### **What is general academic vocabulary?**

In this article, general academic vocabulary is defined as items that have high frequency, wide range, and even distribution in a corpus representing materials from different academic subject areas. These items are identified based on statistical measures (frequency, range, dispersion). Frequency indicates the number of occurrences of a word in the whole academic corpus. The higher frequency a word has, the more likely learners encounter the word in their academic study. However, some words may have high frequency because they are overused in a certain academic subject area, not because they are widely used in a large number of subject areas. For example, photon had a high frequency (600) in the first academic spoken corpus which was used to develop the ASWL, but it appeared in only 6 out of 24 subject areas. Range helps to eliminate these items because it indicates the number of different subject areas in which a word occurs. Yet, range only detects whether the word appears in a subject area or not. It does not discriminate words having different distribution within multiple subject areas. For instance, although predator met the range and frequency criteria of the ASWL, it did not evenly distribute across the first academic spoken corpus. This word occurred 47.38 times per millions in Management, but fewer than 8 times per millions in the remaining 23 subject areas. Dispersion helps to solve this problem because it shows how evenly a word distributes across a corpus. For these reasons, frequency, range, and dispersion have been widely used to identify general academic vocabulary (e.g., Coxhead, 2000; Gardner & Davies, 2014).

Statistical measures allow replicability and comparisons between the lists developed from these studies and those using other corpora or other criteria. They also provide researchers, teachers, and learners with precise and useful information about the occurrences of words in academic discourse. One criticism of using statistical measures to identify items from general academic word lists is that it does not provide the information about variations in the meanings and functions of a word across discourse types. For example, words such as idea, fact, and issue are frequent in both academic and non-academic discourse but may be used differently in each discourse type, or among academic registers. However, general academic word lists take into account that learners might meet these words very often in texts from their specific subject

areas, and knowledge of the core meaning may facilitate the acquisition of discipline-specific meanings. Therefore, while it is important to acknowledge that language use is complex, and it may take some effort to see the relationship between the meaning and use of a word in academic and non-academic discourse, it is equally important not to undermine the value of using statistical measures to distinguish academic and non-academic words.

There are two statistical approaches towards defining academic vocabulary. They reflect two different views towards the relationship between academic vocabulary and general vocabulary (Nation, 2016). One approach (Coxhead, 2000) considers academic vocabulary as part of general vocabulary. According to this view, general vocabulary is a series of layers, each of which represents a 1,000-item frequency band. Words at the 1<sup>st</sup> 1,000-word level are the most frequent and wide ranging items, while those at the 2<sup>nd</sup> 1,000 are less frequent and narrower ranging. The further the 1,000-word levels are from the 1<sup>st</sup> 1,000-word level, the less frequent items in these levels become. Nation (2013) considers items at the 1<sup>st</sup> and 2<sup>nd</sup> 1,000-word levels (e.g., know, sure) as general high-frequency words while Schmitt and Schmitt (2014) argue that these words should be extended to include those at the 3<sup>rd</sup> 1,000-word level. Because general academic wordlists following this approach have set 2,000 as the cut-off point of general high-frequency vocabulary, in this article, we followed Nation's (2013) definition of general high-frequency words. Academic words are defined as items that are outside general high-frequency words but have wide range and high frequency in academic texts. In other words, this approach assumes that learners already know general high-frequency vocabulary and seeks to identify lower frequency words that have wide range and high frequency in academic texts. Many general academic wordlists have been developed using this approach (e.g., Coxhead, 2000; Nesi, 2002).

The second approach considers academic vocabulary as a separate kind of vocabulary that cuts across different 1,000-word levels of general vocabulary (Gardner & Davies, 2014). Academic vocabulary is not seen in the relationship with general high-frequency words. In other words, this approach does not assume that learners know general high-frequency vocabulary. Instead of relying on ready-made lists to distinguish general high-frequency words from general academic words, all items that have wider range and higher frequency in academic rather than non-academic texts are included.

Both approaches provide useful ways of determining the most frequent and wide ranging words in academic texts for EAP learners. The first approach takes into account learners' knowledge of general high-frequency vocabulary and enables learners and teachers to avoid repeatedly learning and teaching known items. In contrast, the second approach allows academic wordlists to avoid limitations related to ready-made general high-frequency wordlists and takes into account the variation in the linguistic features across different discourse types.

These approaches share the same limitation. They look at learners as a homogeneous group that have the same vocabulary knowledge when learning items from their lists. Research has shown that the vocabulary knowledge of L2 learners is diverse. While some learners are able to master at least the most frequent 2,000 words (Laufer, 1998), others have difficulty mastering the most frequent 2,000 words (Henriksen&Danelund, 2015; Matthews & Cheng, 2015; Nguyen & Webb, 2016), and even the most frequent 1,000 words (Henriksen&Danelund, 2015; Nurweni& Read, 1999; Webb & Chang, 2012). Wordlists should suit the level of list users (Nation, 2016). Variation in L2 learners' vocabulary knowledge indicates a need for a general academic wordlist which is adaptable to learners' proficiencies.

### **What wordlists are available to support L2 learners' comprehension of academic spoken English?**

A large number of general academic wordlists (e.g., Coxhead, 2000; Gardner & Davies, 2014) have been developed based on academic written English corpora. There is a large difference in the coverage of the AWL in academic speech (around 4%) (Dang & Webb, 2014; Thompson, 2006) and in academic writing (around 10%), which suggests that lists of academic written words may not be representative of academic spoken vocabulary. Mauranen's (2004) experiment with a highly-experienced oral skills teacher and her EAP class showed that both the teacher and students assumed that items common in written academic text would also be common in academic speech, but then did not find many of these items in an academic spoken corpus. Simpson-Vlach and Ellis (2010) compared their written Academic Formulas List (AFL) and spoken AFL, and found only a 29.07% overlap between them. Lexico-grammar research has also reported a clear-cut difference between the linguistic features of academic speech and academic writing (Biber, 2006; Biber, Conrad, Reppen, Byrd, & Helt, 2002). Taken together, these findings suggest that a wordlist which is developed from a written corpus may not capture



the language in academic speech as fully as a wordlist developed from an academic spoken corpus.

Despite this fact, little effort has been made to develop an academic spoken wordlist. To the best of our knowledge, only two studies have focused on creating academic spoken wordlists. Nesi (2002) developed a Spoken Academic Word List (SAWL) of single-words, but unfortunately, to date, there are no descriptions of the development, validation or items in her list. Simpson-Vlach and Ellis (2010) focused on multi-word units by creating a spoken AFL. This has great value because knowledge of multi-words is essential for fluent processing (Simpson-Vlach & Ellis, 2010). Yet, knowledge of single-words is also important, because it provides valuable support for the acquisition of multi-words. Although phrases in different lists of multi-words may vary, they share a considerable number of core single-words (Adolphs & Carter, 2013; Shin & Nation, 2008). Therefore, it is beneficial to create an academic spoken list of single-words.

### **The scope of the present study**

The ASWL presented in this paper is a general academic wordlist. It is aimed towards (1) EGAP programs, (2) EAP programs where learners are unclear about their target subject areas, or teachers lack background knowledge of learners' specific disciplines, and (3) interdisciplinary environments where it is unclear which specific discipline an academic subject belongs to. The development of the list expands on the two statistical approaches towards identifying academic vocabulary. It views general academic vocabulary as a separate kind of vocabulary that cuts across various frequency levels of general vocabulary, and therefore, develops the list from scratch. It also considers academic spoken vocabulary in relation to general vocabulary by making the list adaptable to learners' knowledge of general vocabulary.

### **Research questions**

1. Which lexical items occur frequently and are evenly distributed in a wide range of academic speech?
2. What is the coverage of the ASWL in independent collections of academic speech, academic writing, and non-academic speech?
3. With knowledge of the ASWL, how much coverage of academic speech may be reached by learners with different levels of general vocabulary?

## **Methodology**

### **Corpora development**

Four corpora of around the same size were compiled in this study: two academic spoken corpora, one academic written corpus, and one non-academic spoken corpus. The first academic spoken corpus was used to create the ASWL while the other corpora were used to validate the list from different perspectives (see Table 1). This satisfies Nation and Webb's (2011) guideline that a list should be validated in an independent corpus of similar size as the corpus from which it was developed.

[TABLE 1 NEAR HERE]

The ASWL aims to help EAP learners from different academic disciplines to enhance their comprehension of academic speech in English-medium university programs. Therefore, the academic spoken corpora should represent speech events from a wide range of academic disciplines that these students are likely to encounter often in their future study. To achieve this goal, materials in the two academic spoken corpora were selected from 11 sources which represent naturally occurring academic speech recorded in various institutions around the world and represent a wide range of varieties of English (see Table 1 (Supporting Information online)). These materials were written transcripts of spoken data recorded by other researchers rather than the current researchers themselves to deal with the challenge of developing academic spoken corpora (McCarthy & Carter, 1997; Thompson, 2006). Whole texts rather than partial texts were included because samples of whole texts better reflect the target language than partial texts due to the variation in the linguistic features across different parts of the text (Biber, 2006; Sinclair, 1991).

The two corpora have very similar sizes and structures so the validating corpus reflects closely the vocabulary in the corpus from which the list is developed, and provide an accurate assessment of the ASWL. Each corpus contains about 13-million running-words; meaning they are eight times larger than Nesi's (2002) SAWL corpus (1.6-million), more than six times larger than Simpson-Vlach and Ellis's (2010) spoken AFL corpus (2.1-million), and four times larger than Coxhead's (2000) AWL corpus (3.5-million). Given the wide recognition of the AFL and

AWL, it is expected that the two academic spoken corpora in this study are large enough to capture the most frequent and wide ranging words in academic spoken English.

In terms of representativeness, Biber's (1993) and Coxhead's (2000) guidelines were followed so that the two academic spoken corpora represent as closely as possible the academic speech that EAP learners from a wide range of academic disciplines are likely to encounter in their academic study in English-medium programs. Each corpus has two levels and represents four kinds of speech events.

At the macro level, the corpus is divided into four disciplinary sub-corpora based on Becher's (1989) classification of academic disciplines in higher education: hard-pure, hard-applied, soft-pure, and soft-applied. The hard/soft dimension is related to the degree to which a paradigm exists. The pure/applied dimension is associated with application to practical problems. Becher's classification has been validated in a wide range of contexts (Jones, 2011). Adopting this classification to construct the two academic spoken corpora ensures that their structure is not biased towards the administrative structure of a particular institution. This broad level allows the identification of words common between disciplines and compare their occurrences in each discipline. At the micro level, each disciplinary sub-corpus is divided into a number of subject areas to ensure a wide range of academic subjects. Ideally, each disciplinary sub-corpus should consist of the same number of equally-sized subject areas so that the ASWL will not be biased toward the vocabulary in a certain discipline or subject area.

Four kinds of speech events are represented: lectures, seminars, labs, and tutorials. Lectures are the most common academic speech events and are opportunities in which lecturers inform, evaluate, and critique important information in the reading materials that they would like to draw their students' attention to (Lynch, 2011). In this study, lectures are defined as events in which lecturers are the ones who mainly speak. Seminars, tutorials, and labs are opportunities for students to participate in group discussion with lecturers, tutors, and fellow students (Adolphs & Carter, 2013). The target users of the ASWL are L2 learners rather than experts. Therefore, in this study, seminars are defined as student instructional seminars, involving interactions between course instructors and students where students are the ones who mainly speak (Aguilar, 2016) rather than expert research seminars, which are opportunities for academics to speak about their on-going or completed research to a small expert audience

(Aguilar, 2016). Labs and tutorials are the speech events that provide students with opportunities to deepen their understanding of the information from the lectures and develop practical skills. Labs are more common in hard subjects; tutorials are more common in soft subjects (Neumann, 2001). Including these four speech events ensures that the two academic spoken corpora represent both common speech events across disciplines and distinctive speech events in each discipline. Lectures and seminars respectively account for the largest and second largest proportion in each sub-corpus of the academic spoken corpora. Next come labs and tutorials (depending on whether the disciplines are hard or soft). This proportion is aligned with the proportion of these speech events in the six corpora from publishers and that stated in the 2013 undergraduate and postgraduate course outlines at the current researchers' institution.

Lectures, seminars, labs, and tutorials rather than outside classroom speech events (e.g., office hours, service encounters) were chosen to represent academic spoken English because they address the primary need of the ASWL target users to be able to comprehend and engage in the academic courses with their instructors and fellow students. Given the limited time and slow vocabulary growth rates of EFL learners (Milton, 2009; Webb & Chang, 2012), focusing on the vocabulary in classroom speech events has practical value. It also allows the classification of the texts into the four disciplinary sub-corpora because classroom language is more subject-focused than outside classroom language (Csomay, 2006). To maximize the representativeness of the two academic spoken corpora, when possible, materials from all 11 sources and different varieties of English were presented in each sub-corpus of the two corpora.

Tables 2 and 3 (Supporting Information online) present the composition of the first academic spoken corpus (to develop the ASWL) in terms of disciplines and speech events, respectively. This corpus consists of four sub-disciplinary sub-corpora. Each sub-corpus has around 3.25-million running-words from about 380 texts, and is divided into six subjects. Each subject contains around 500,000 running-words. Six subjects per sub-corpus is sufficient. It is around the same or even larger as the number of subjects per sub-corpus in the academic corpora of previous research (Biber, 2006; Coxhead, 2000; Hyland & Tse, 2007). The texts in the first academic spoken corpus come from all 11 sources and represent at least seven varieties of English. All four speech events are represented. Lectures have the largest proportion of texts. Next come seminars. Labs and tutorials have the smallest proportion.

Tables 4 and 5 (Supporting Information online) demonstrate the composition of the second academic spoken corpus which was used to validate the ASWL. Similar to the first academic spoken corpus, this corpus has around 13-million running-words and is divided into four disciplinary sub-corpora. Each contains around 3.2-million running-words. This corpus was less balanced than the first academic spoken corpus in terms of the number of subjects per sub-corpus and the number of words per subject. All four kinds of speech events were represented with lectures and seminars making up larger proportions of texts than labs and tutorials.

Tables 6-8 (Supporting Information online) presents the structure of the academic written corpus and the non-academic spoken corpus. The academic written corpus represents different kinds of academic writing (book chapters, journal articles, student writings, research reports, and textbooks) in courses at four different locations. It has similar structure as the two academic spoken corpora with four disciplinary sub-corpora. Each contains more than 3-million running-words and represents a range of subject areas. The non-academic spoken corpus was developed to examine if the ASWL reflects academic vocabulary. It is comprised from seven sources which represent different kinds of general spoken English and 10 varieties of English.

### **Determining the unit of counting for the ASWL**

Two common units of counting in academic wordlists are lemmas and word-families. A lemma (predict) consists of a stem (predict) together with its inflected forms (predicts, predicted, predicting). Members of a lemma belong to the same word class (Francis & Kučera, 1982). A word-family (predict) consists of a stem (predict), its inflections (predicts, predicted, predicting), and closely related derivations up to Level-6 of Bauer and Nation's (1993) scale (predictably, predictable, unpredictable, unpredictably, predictability, prediction, predictions, predictive, predictor, predictors, unpredictability). Word-families were the unit of counting of many earlier general academic wordlists (Coxhead, 2000; Nesi, 2002; Xue & Nation, 1984). It is also the common unit of counting in numerous general vocabulary lists and discipline-specific wordlists (e.g., Liu & Han, 2015; Nation, 2012; Wang, Liang, & Ge, 2008; Yang, 2015). Researchers (Brezina & Gablasova, 2015; Gardner & Davies, 2014; Lei & Liu, 2016) have recently questioned the suitability of word-families as a unit of counting and proposed using lemmas instead. Each of their criticisms has been addressed by Nation (2016) in detail.

The first criticism is that word-families are not as appropriate a unit of counting as lemmas because not all members of a word-family are closely related in meaning (Gardner & Davies, 2014). Nation (2016) points out that, to be included in a word-family in Bauer and Nation's (1993) scale, the meaning of the base in the derived word must be closely related to the meaning of the base when it stands alone or is combined with other derived forms.

The second criticism is that word-families do not make part of speech distinctions while lemmas do (Gardner & Davies, 2014; Lei & Liu, 2016). Nation (2016) argues that distinguishing part of speech has a negative impact on the distinction of very closely related items like walk (v) and walk (n), but cannot distinguish homonyms having the same part of speech like bank (for money) and bank (for river).

The third criticism is that learners, even young native speakers, may not have knowledge of word building devices of English, and therefore, lemmas is a more suitable unit of counting for them (Brezina & Gablasova, 2015; Gardner & Davies, 2014). Nation (2016) points out that the word-family in Bauer and Nation's scale is a set of levels, which is based on frequency, productivity, predictability, and regularity of affixes. In this scale, word-families are divided into seven levels with Level 1 consisting of single word types with no family members, Level-2 consisting of the most elementary and transparent members of a word-family, and Level-7 consisting of the least transparent members. The lemma is Level-2 word-family while the word-family referred to in previous studies (e.g., Coxhead, 2000) is Level-6 word-families.

[TABLE 2 NEAR HERE]

Table 2 presents the list of affixes at each word-family level. If a word-family is defined as being at a particular level, it will include the stem together with its potential inflections and derivations made up of affixes up to that level. For example, members of a Level-6 word-family would be the stem itself, and can potentially include members derived from one or more affixes up to Level-6 (eight from Level-2, 10 from Level-3, 11 from Level-4, 50 from Level-5, and 12 from Level-6). In other words, lemmas (Level-2 word-families) and Level-6 word-families reflect different steps toward the full morphological knowledge (Level-7). Therefore, the question is not whether lemmas or word-families are better, but which word-family level is the most suitable for a particular group of learners.

To find the answer to this question, the learning burden (i.e., the amount of effort needed to acquire a word-family at a certain level) should be taken into account (Nation, 2013, 2016; Nation & Webb, 2011). The idea behind word-families is that learners may be able to see that word forms with the same stems are related to each other, and therefore, may find it easier to recognize or learn a word which is morphologically related to a known word rather than a totally unrelated word. Therefore, while it is important to recognize that it takes some effort to see the relationships between some members of the word-families, it is equally important not to overstate the differences. According to Nation (2016), it is not difficult to infer the meaning of walk (n) in When I go for a walk if the meaning of walk (v) is known. Likewise, minimal efforts are needed to infer the meaning of sadness and sadly if learners already know sad and -ly and -ness, and have met these affixes in several words.

Word-families up to Bauer and Nation's (1993) Level-6 were chosen as the primary unit of counting for the ASWL for three reasons. First, following Coxhead (2000) and Nation (2013), the use of word-families at Level-6 in this study is looked at from a pedagogical perspective. That is, knowledge of word-family members is gradually picked up during the learning process rather than acquired all at the same time, and learners are provided with training on word part knowledge and word building skills. In this way, knowledge of one word-family member will help learners to facilitate the acquisition of other members. This assumption is supported by earlier studies showing that L2 learners' derivational knowledge increased incrementally over time (Mochizuki & Aizawa, 2000; Schmitt & Meara, 1997; Schmitt & Zimmerman, 2002), and instructions of word parts helped to expand learners' vocabulary knowledge (Schmitt & Meara, 1997; Wei, 2014).

Second, one purpose of this study is to integrate the ASWL with Nation's (2012) BNC/COCA lists to organize a more systematic language program to support L2 learners' comprehension of academic spoken English. The Level-6 word-family is the unit of counting of the BNC/COCA2000. It has been widely used as the unit of counting in vocabulary research and the design of learning and teaching materials for L2 learners. Choosing the Level-6 word-family as the unit of counting of the ASWL allows learners and teachers to make good use of numerous available resources.

Third, research on derivational knowledge of L2 learners from different L1 backgrounds and learning contexts has shown that L2 learners, even beginners, do know a number of affixes at Levels 3-6 in Bauer and Nation's (1993) scale (Mochizuki & Aizawa, 2000; Sasao & Webb, 2017). For example, re-(Level-6), -ful(Level 4), and un-(Level 5) were known by 70%-75% of the beginner learners in Mochizuki and Aizawa's (2000) study. These affixes were also categorized as the beginner level in terms of difficulty by Sasao and Webb (2017) who measured the word part knowledge of 1,348 people representing more than 30 different L1s. The ASWL aims to benefit learners with different proficiency levels. Choosing Level-2 word-families (lemmas) may then overestimate the learning burden of a word-family for most ASWL target users. One question that arises is which level from Level 3 to Level-6 is the most suitable unit of counting for the ASWL. Bauer and Nation's (1993) scale was based on usefulness and regularity, not learner knowledge. Research has shown that L2 learners' knowledge of affixes does not neatly fit in this scale but varies according to learners' L1 and L2 proficiency, and the instruction of word parts that they have received (Mochizuki & Aizawa, 2000; Sasao & Webb, 2017). Therefore, while Bauer and Nation's (1993) scale is a useful framework, it should be applied with flexibility when creating pedagogical wordlists (Nation, 2016). The ASWL target users are from different learning contexts. Choosing Level-6, which is nearly the broadest word-family level, as the unit of counting of the ASWL may deal with the diversity in the characteristics of the list users to some extent.

While the Level-6 word-family was chosen as the primary unit of counting of the ASWL, we acknowledge that this unit of counting may overestimate the morphological knowledge of a proportion of EFL learners (Mochizuki & Aizawa, 2000; Schmitt & Meara, 1997; Schmitt & Zimmerman, 2002; Ward & Chuenjundaeng, 2009). Thus, apart from the list of Level-6 word-families which is presented in this article, another version of the ASWL which lists the Level-2.5 word-families was also created. Both versions of the ASWL will be freely available at our websites. The Level-2.5 word-family is relevant to the lemma but does not distinguish part of speech (Nation, 2016). Pinchbeck (2014) calls this unit of counting flemma with f standing for family. For instance, form (verb) and form (noun) are counted as two lemmas. However, they are considered as one flemma. Flemmas rather than lemmas were chosen because, as mentioned, distinguishing part of speech may overestimate the learning burden of very closely related items like walk (v) and walk (n) but cannot distinguish homonyms having the same part of speech such



as bank (for money) and bank (for river). Moreover, the use of flemmas also makes it possible to do the analysis with Nation, Heatley, and Coxhead's (2002) RANGE program. The flemmas list was created by grouping the ASWL Level-6 word-family members. Leech, Rayson, and Wilson's (2001) principles for creating lemmatized wordlists were used as a guide. For example, the Level-6 word-family acquire has seven members: acquire, acquired, acquires, acquiring, acquirer, acquirers, and unacquired. When converted into flemmas, they were grouped into three Level-2.5 word-families: acquire (acquire, acquired, acquires, acquiring), acquirer (acquirer, acquirers), and unacquired (unacquired).

### **Key considerations for the ASWL**

To ensure that the ASWL can benefit a wide range of EAP learners planning to study different academic subjects and having different proficiencies, the list must have four following characteristics.

- (1) Size and coverage: The ASWL must contain a smaller number of word-families but provide higher coverage in the first academic spoken corpus than Nation's (2012) BNC/COCA2000. The BNC/COCA2000 contains items from the 1<sup>st</sup> and 2<sup>nd</sup> 1,000 BNC/COCA frequency levels. To ensure that the list represents the vocabulary that L2 learners encounter, Nation (2012) developed the BNC/COCA2000 from a 10-million running-word corpus, 60% spoken (spoken English, movies, and TV programs), and 40% written (texts for young children and fiction). The corpus represents three varieties of English: American-English, British-English, and New Zealand-English. Dang and Webb (2016) and Dang (2017) compared the BNC/COCA2000 with three other general high-frequency wordlists (West's (1953) General Service List, Nation's (2004) BNC2000, and Brezina and Gablasova's (2015) New-GSL) using lexical coverage, teacher perceptions of word usefulness, and learner vocabulary knowledge as criteria. The results suggested that the BNC/COCA2000 was the most suitable general high-frequency wordlist for L2 learners; therefore, it was chosen to represent general high-frequency vocabulary in the present study.
- (2) Word-families outside general high-frequency word-families: The ASWL must include a considerable number of word-families outside the BNC/COCA2000, but have high frequency, wide range, and even distribution in the first academic spoken corpus.

Criteria 1 and 2 were established because the ASWL aims to direct EAP learners' attention to the most frequent and wide ranging words in academic speech that are beyond their existing vocabulary levels. Therefore, if the ASWL either had a larger size but provided lower coverage than the BNC/COCA2000 or contained mainly the BNC/COCA2000 words, it would not draw much interest from EAP learners and teachers. They may simply use the BNC/COCA2000 rather than putting their effort into a new list.

- (3) Distribution across the four sub-corpora: The coverage of the ASWL in the sub-corpora of the first academic spoken corpus should be similar. In this way, the ASWL can benefit a wide range of EAP learners irrespective of their disciplines.
- (4) Adaptability to *learners' levels*: The ASWL must be divided into four levels according to Nation's (2012) BNC/COCA lists to benefit learners with different proficiencies. Levels 1-3 contain ASWL word-families from the 1<sup>st</sup>, 2<sup>nd</sup>, and 3<sup>rd</sup> 1,000 BNC/COCA frequency levels, respectively. Level-4 represents ASWL items that are outside the most frequent 3,000 BNC/COCA word-families. To make sure that the list benefits learners at different vocabulary levels, not only the whole list but its levels should have the first three characteristics.

### **Developing and validating the ASWL**

Nation et al.'s (2002) RANGE was used to count and sort the words in the first academic spoken corpus. This program is available at Paul Nation's website (<http://www.victoria.ac.nz/lals/about/staff/paul-nation>). Because the present study aims to make a corpus-based wordlist that is suitable for L2 learning and teaching, Nation's (2016) principles were followed in the ASWL development. That is, objective criteria from corpora should be used as the main criteria in word selection; however, these criteria should be adjusted based on human judgement to make the list more useful and suitable for the target users.

The development of the ASWL is primarily based on the analysis of the range, frequency, and dispersion of the word-families in the first academic spoken corpus. These criteria allow the ranking of all the word-families in the first academic spoken corpus from the most frequent and wide ranging items to the least frequent and narrowest items. Ideally, learners should gradually learn from the most frequent and wide ranging words to the least frequent and narrow ranging words so that they can acquire items that they are likely to encounter often in their academic

study before those they are least likely to encounter. However, it will be impractical to learn all items in the corpus given that L2 learners have less learning time and slower vocabulary growth rates than L1 children (Milton, 2009; Webb & Chang, 2012). Therefore, different pilot versions adopting different range, frequency, and dispersion cut-off points were compared. The four key considerations for the ASWL (size and coverage, word-families outside general high-frequency word-families, distribution across the four sub-corpora, and *adaptability to learners' levels*) were used as the guide to see which version would have the greatest pedagogical value.

### **Range**

A selected word-family had to occur in all four disciplinary sub-corpora (hard-pure, hard-applied, soft-pure, and soft-applied) of the first academic spoken corpus, and at least 50% of the subject areas (12 out of 24 subjects). This criterion has been consistently used in studies aimed at developing specialized wordlists (e.g., Coxhead, 2000; Coxhead & Hirsh, 2007; Wang et al., 2008). It ensures that students from a wide range of disciplines and subject areas can gain benefit from learning these lists.

### **Frequency**

A selected word-family had to occur at least 350 times in the first academic spoken corpus (at least 26.9 times per million words). This frequency figure was the result of extensive experimentation which compared the items included or excluded from the ASWL at different frequency cut-off points from 50 to 370. The 370 figure was transferred from the frequency criterion used by previous studies to select items for academic written wordlists (e.g., Coxhead, 2000).

First, the pilot version with the frequency cut-off point of 350 better fulfilled the first characteristic of an ASWL (size and coverage) than the versions with the frequency cut-off points of 250 or lower. Although all of them provided higher coverage than the BNC/COCA2000, the version with the frequency cut-off point of 350 had fewer than 2,000 items (1,741 word-families) while the other versions had more than 2,000 items (2,033-3,416 word-families).

Second, the pilot version with the frequency cut-off point of 350 satisfied the fourth characteristic of an ASWL (*Adaptability to learners' levels*) better than the frequency cut-off point of 300. The former version is more beneficial for learners who have mastered the most

frequent 1,000 words than the latter version. Although both versions provided higher coverage than the 2<sup>nd</sup> 1,000 BNC/COCA word-families, the version with the frequency cut-off point of 350 had fewer items beyond the 1<sup>st</sup> 1,000 BNC/COCA frequency level (911 word-families) than the 2<sup>nd</sup> 1,000 BNC/COCA word-families while the list with the frequency cut-off point of 300 had more items (1,024 word-families).

Third, compared with the frequency cut-off point of 370, 350 offered a better compromise between the second characteristic (Word-families outside general high-frequency word-families) and third characteristic (Distribution across the four sub-corpora). The examination of the distribution of these versions (both as a whole and each level) showed that the version with the frequency cut-off point of 350 satisfied the third characteristic (Distribution across the four sub-corpora) as well as the version with the frequency cut-off point of 370. However, lowering the frequency cut-off point from 370 to 350 means adding 23 word-families outside the most frequent 2,000 BNC/COCA word-families to the ASWL (e.g., perception, infer, analytic, subtract). These word-families came from the 3<sup>rd</sup> to 7<sup>th</sup> 1,000-BNC/COCA word levels, and most of them occurred in at least 20 out of 24 subjects and had a dispersion of 0.7 or above.

### **Dispersion**

A selected word-family had to have Juilland and Chang-Rodrigues's (1964) dispersion of at least 0.6 across 24 subjects. Although there are several ways to examine dispersion, Juilland and Chang-Rodrigues's (1964) dispersion is the most common way to measure dispersion in studies developing specialized wordlists (e.g., Gardner & Davies, 2014; Lei & Liu, 2016) and general wordlists (e.g., Nation, 2006; Nation, 2012). The value of Juilland and Chang-Rodrigues's (1964) dispersion can range from 0 (extremely uneven distribution) to 1 (perfectly even distribution). Similar to the frequency criterion, the dispersion cut-off point of 0.6 was chosen based on extensive comparison of different pilot versions with different dispersion cut-off points from 0.1 to 0.9.

A pilot version with the dispersion cut-off point of 0.6 fulfilled the first characteristic of an ASWL (Size and coverage) better than those with the dispersion cut-off points of 0.1 and 0.2. Although all cut-off points resulted in pilot versions that provided higher coverage than the BNC/COCA2000, the version with the dispersion cut-off point of 0.6 contained fewer than 2,000

items (1,741 word-families), but those with the dispersion cut-off points of 0.1 and 0.2 had more than 2,000 items (2,023 word-families, 2,016 word-families).

The 0.6 cut-off point better satisfied the third characteristic (Distribution across the four sub-corpora) and the fourth characteristic (*Adaptability to learners' levels*) than the cut-off points of 0.3, 0.4, 0.5, and 0.7. Comparison of the coverage of these versions (both as a whole and each level) in the four sub-corpora showed that the cut-off point of 0.6 either ranked first or second in terms of the degree of evenness in distribution across the sub-corpora. In contrast, the other versions had lower or unstable ranking, or both.

Third, the cut-off point of 0.6 met the second characteristic (Word-families outside general high-frequency word-families) better than the pilot versions with the dispersion cut-off point of 0.8 and of 0.9. The cut-off point of 0.6 had 455 words families outside the BNC/COCA2000, and their coverage in the first academic spoken corpus was 3.28% (whole corpus) and 2.63%-3.64% (sub-corpora). These word-families represented items from the 3<sup>rd</sup> to 9<sup>th</sup> 1,000 BNC/COCA word levels, and one word-family outside the most frequent 25,000 BNC/COCA word-families. In contrast, the cut-off point of 0.8 had a much smaller number of word-families outside the BNC/COCA2000 (154). These word-families appeared at the 3<sup>rd</sup>, 4<sup>th</sup> and 10<sup>th</sup> 1,000 BNC/COCA word levels. These word-families covered only 1.19% of the 1<sup>st</sup> academic spoken corpus and 1.16%-1.25% of its sub-corpora. Similarly, 99.38% of the items in the pilot version with the dispersion cut-off point of 0.9 were from the BNC/COCA2000.

In sum, the frequency cut-off point of 350 and the dispersion cut-off point of 0.6 are the best compromise between the four characteristics of the ASWL. These cut-off points were lower than those used to select Coxhead's (2000) AWL (frequency of 370) and Gardner and Davies's (2014) AVL words (dispersion of 0.8), which supports the findings of previous research that there is a clear-cut difference between the linguistic features of academic speech and academic writing (Biber, 2006; Biber et al., 2002).

Items that satisfied the range, frequency, and dispersion criteria were included in the ASWL. They were divided into four levels according to the BNC/COCA frequency levels so that the list is suitable to learners at different proficiency levels. Each level was divided into a sub-list of function words, and sub-lists of lexical words. The lexical words were further divided into sub-lists of 50 items according to the frequency of the word-families in the first academic spoken

corpus. The coverage of the ASWL and its levels in the two academic spoken corpora, the academic written corpus, and the non-academic spoken corpus was determined by running these corpora in turn through RANGE with the ASWL and its levels as the base wordlists.

## Results

### **RQ1. Which lexical items occur frequently and are evenly distributed in a wide range of academic speech?**

In the first academic spoken corpus, 1,741 word-families satisfied the criteria for inclusion in the ASWL. Table 9(Supporting Information online)presents the ASWL headwords at each level. Tables 10-13(Supporting Information online)present headwords in each sub-list within each level. The headwords are at Bauer and Nation’s (1993) Level-6. All ASWL word-families occur in at least 14 out of the 24 subjects in the corpus. 75.36% of the ASWL occur in all 24 subjects, and 97.47% occur in 20 or more. In terms of dispersion, although 0.6 was set as the cut-off point, 85.41% of the ASWL words have a dispersion of at least 0.7.

Table 3 presents the lexical profile of the ASWL and its four levels. Levels 1 and 2 contain general high-frequency words from the most frequent 2,000 BNC/COCA word-families that met the range, frequency, and dispersion criteria. Levels 3 and 4 are academic words that have high frequency in academic speech and are outside general high-frequency words.

[TABLE 3 NEAR HERE]

The ASWL accounted for 90.13% of the first academic spoken corpus, which is higher than the coverage of the most frequent 2,000 BNC/COCA word-families (88.61%). It should be noted that the former list has 259 fewer items than the latter list. When tested in each disciplinary sub-corpus, the ASWL consistently provided around 90% coverage: 89.46% (hard-pure), 91.07% (hard-applied), 89% (soft-pure), and 90.92% (soft-applied).

Readers may be interested to know how the ASWL is compared with Coxhead’s (2000) AWL and Gardner and Davies’s (2014) AVL. Following Gardner and Davies(2014), the Level-6 word-family version of the AVL was used for the comparison. This version was downloaded from <http://www.academicvocabulary.info/>. In this version, Gardner and Davies distinguish between word classes. However, in this study, to be consistent, repeated word-families were removed, and the final version of the AVL contained 1,983 word-families. Table 4 presents the coverage of

these lists in the two academic spoken corpora. The ASWL provided higher coverage in the two academic spoken corpora (around 90%) than the AVL (around 4%). Moreover, although the ASWL had 242 fewer items than the AVL, it provided higher coverage in the two academic spoken corpora than the AVL (around 24%). These findings may be due to either the difference between academic spoken and written English or the difference in the principles behind the development of the three lists.

[TABLE 4 NEAR HERE]

**RQ2. What is the coverage of the ASWL in independent collections of academic speech, academic writing, and non-academic speech?**

The ASWL covered around 89.59% of the words in the second academic spoken corpus, which is similar to its coverage in the first academic spoken corpus (90.13%). In contrast, the ASWL provided lower coverage in the academic written corpus (81.43%) and the non-academic corpus (87.06%). This findings is consistent with the findings of previous research on developing academic wordlists (Coxhead, 2000; Gardner & Davies, 2014). That is, the coverage of the list in the corpus from which it was developed is similar to its coverage in an independent corpus representing the same types of discourse, but higher than its coverage in independent corpora representing different kinds of discourse.

**RQ3. With knowledge of the ASWL words, how much coverage of academic speech may be reached by learners with different levels of general vocabulary?**

Table 5 demonstrates the potential coverage that learners with different vocabulary levels may achieve with the aid of the ASWL. This potential coverage is the sum of the coverage provided by two groups of words. The first group includes the word-families that students may already know. They are items in the BNC/COCA levels that are at students' existing vocabulary level. The second group is the coverage provided by the ASWL word-families that students may not know. They are ASWL items that are outside the BNC/COCA words in the first group. The second column of Table 5 presents the number of ASWL word-families that are beyond learners' existing vocabulary levels. The next two columns show the potential coverage that learners may reach if they learn these ASWL words. Coverage provided by proper nouns (e.g., John, Berlin) and marginal words (e.g., ah, hmm) is presented in the last two rows. Previous research on the vocabulary load of spoken English (e.g., Nation, 2006) assumed that proper nouns and marginal

words have a minimal learning burden for learners, and therefore, added the coverage by these words to the potential coverage. The last two columns demonstrate the potential coverage including proper nouns and marginal words.

[TABLE 5 NEAR HERE]

The first row of Table 5 shows the potential coverage for learners who are yet to master the most frequent 1,000 BNC/COCA word-families. These learners are not likely to know the ASWL due to their insufficient vocabulary knowledge. If they study all 1,741 ASWL word-families, they may reach 90% coverage of the two academic spoken corpora. If proper nouns and marginal words are known, the potential coverage for these learners is 92%-93%. These coverage figures are higher than those provided by the most frequent 2,000 BNC/COCA word-families (91.08%, 90.54%).

The potential coverage for learners who have mastered the most frequent 1,000 BNC/COCA word-families is presented in the second row of Table 5. With their existing knowledge, these learners only need to learn 911 ASWL word-families. Yet, they may gain potential coverage of 90%-91%. Including proper nouns and marginal words, the potential coverage that these learners may reach with the aid of the ASWL is 93%. It is higher than the potential coverage if they studied 1,000 word-families from the 2<sup>nd</sup> 1,000 BNC/COCA frequency level instead. Interestingly, learning the ASWL still allows these two groups of low-level learners to gain 92.21% (those with the vocabulary level less than 1,000 word-families) and 93.45% (those with the vocabulary level of the most frequent 1,000 word-families) coverage of non-academic speech.

The third row of Table 5 shows that, learners with knowledge of the most frequent 2,000 BNC/COCA word-families only have to learn 455 words in the ASWL, but can gain the potential coverage of 91%-92% (without proper nouns and marginal words) and 94% (with proper nouns and marginal words) of academic speech. For those with knowledge of the most frequent 3,000 BNC/COCA word-families (see the fourth column of Table 5), knowledge of 75 ASWL word-families that are beyond their vocabulary levels may enable them to achieve 92%-93% coverage of academic spoken discourse. If proper nouns and marginal words are counted, the potential coverage is 95%-96%. Taken as a whole, the potential coverage including proper nouns and marginal words ranges from 92% to 96%.



## **Discussion**

### **Is there a core vocabulary in academic spoken English?**

The ASWL covered around 90% of the words in the corpus from which it was developed and an independent academic spoken corpus of a similar size and structure. This finding indicates that the ASWL accurately captures high-frequency, wide-ranging, and evenly-distributed word-families in academic speech. The higher coverage of the ASWL in the two academic corpora than in the academic written corpus suggests that the list better represents spoken than written vocabulary. Similarly, its coverage in the two academic spoken corpora was higher than in the non-academic spoken corpus. This demonstrates that the ASWL accurately represents academic rather than non-academic vocabulary. The ASWL was developed from an analysis of academic speech from a wide range of subject areas but still provided similar coverage in hard-pure, hard-applied, soft-pure, and soft-applied subjects. This finding indicates that there is a core vocabulary across academic speech of different disciplines, and the list can offer fairly equal benefit to EAP learners regardless of their subject areas.

Moreover, the disciplinary division of the two academic spoken corpora used to develop and validate the ASWL is based on Becher's (1989) classification of academic disciplines in higher education. The validity of Becher's classification has been confirmed in various contexts, which indicates that it is transferable across institutions and can serve as a common standard for comparison (Nesi, 2002). As a result, this classification has been widely used as a way to categorize academic disciplines in higher education (see Jones, 2011 for more details) and to structure academic corpora such as the BASE and British Academic Written English (BAWE) corpora and Hyland's (2000) academic written corpus. Given the high validity and wide transferability of Becher's classification, it is expected that the ASWL can be globally used by EAP learners irrespective of the administrative structure of their universities. In brief, the fact that the ASWL can benefit L2 learners irrespective of their disciplines and institutional structures highlights the value of general academic wordlists for EGAP programs.

### **Can the ASWL benefit learners with different vocabulary levels?**

This study suggests that learners with different vocabulary levels can benefit from the ASWL. Intermediate-level learners may achieve around 95% coverage of academic speech by learning a small number of items from the ASWL: 455 word-families (those with the vocabulary level of

the most frequent 2,000 word-families) or 75 word-families (those with vocabulary level of the most frequent 3,000 word-families). It is important for learners to know at least 95% of the words in spoken discourse (Schmitt, Cobb, Horst, & Schmitt, 2015). It allows them to obtain a high and stable degree of listening comprehension (van Zeeland & Schmitt, 2013). Dang and Webb (2014) found that a vocabulary size of 4,000 word-families is needed to reach 95% coverage of academic speech. This means that learners with knowledge of the most frequent 2,000 word-families may need to learn a further of 2,000 word-families from the 3<sup>rd</sup> and 4<sup>th</sup> 1,000 word levels. Meanwhile, those with knowledge of the most frequent 3,000 word-families may need to study an extra 1,000 word-families from the 4<sup>th</sup> 1,000 word level. Studying the ASWL word-families that are beyond their levels gives these learners a better return. They need to learn a much smaller number of items but are still able to achieve 95% coverage of academic spoken discourse.

Ideally learners would study the most frequent 2,000 or even 3,000 BNC/COCA word-families and then move to the relevant ASWL levels so that they can reach 95% of academic spoken discourse. However, this may be too demanding a goal for low-level learners, especially those studying in EFL contexts. L2 learners may learn an average of 400 word-families (Webb & Chang, 2012) a year. This means that low-level learners may need about six years to acquire the most frequent 2,000 word-families plus 455 extra ASWL word-families, or eight years to acquire the most frequent 3,000 word-families plus 75 extra ASWL word-families. Research with learners in a wide range of EFL contexts such as China (Matthews & Cheng, 2015), Denmark (Henriksen & Danelund, 2015), Indonesia (Nurweni & Read, 1999), Israel (Laufer, 1998), Taiwan (Webb & Chang, 2012) and Vietnam (Nguyen & Webb, 2016) suggests that some learners may have even slower vocabulary growth rates. A reasonable proportion of learners in these studies had not mastered the most frequent 2,000 word-families let alone the most frequent 1,000 word-families after six years or more of formal English instruction. This slow vocabulary growth rate means it may often be challenging for EFL learners to master the 2,455 or 3,075 word-families by the time their academic programs start. Going straight to the ASWL or learning the ASWL from Level-2 can help with this dilemma to some degree. It focuses low-level learners' attention on the most important words in academic speech, meanwhile, allows them to make up for their insufficient knowledge of general high-frequency words.

Learning the ASWL words can help beginners to achieve 90%-91% and 92%-93% coverage of academic speech without and with knowledge of proper nouns and marginal words, respectively. These figures are meaningful for low-level learners in three ways. First, these learners have to study a much smaller number of items, but can achieve higher coverage of academic speech than learning words from the subsequent levels of general vocabulary. Second, although comprehension may not be as easy as with 95% coverage, 90%-93% coverage may still allow L2 learners to achieve basic comprehension of academic speech. Van Zeeland and Schmitt (2013) found no significant difference in L2 listening comprehension between the 90% and 95% coverage figures. Both coverage figures led to 70% listening comprehension. Moreover, if 70% comprehension was considered as an indicator of adequate comprehension, both figures resulted in 75% of the participants achieving this result. Even with a stricter criterion of adequate comprehension (80%), more than half of the participants met this requirement at both coverage levels. Furthermore, in real life academic speech, students receive support to facilitate their listening comprehension such as pre-lecture reading materials, visual aids, or interaction with their lectures, tutors, and peers, which may allow L2 learners to compensate for their inadequate vocabulary knowledge and enhance their listening of academic speech (Flowerdew & Miller 1992; MacDonald, Badger, & White 2000; Mulligan & Kirkpatrick, 2000). Third, going straight to the ASWL or learning the ASWL from Level-2 may also enable beginners to reach 92%-93% of general spoken English, which may allow them to understand this important discourse type. Taken together, the ASWL is a good shortcut for low-level learners to achieve basic comprehension of academic speech while still allowing them to enhance their knowledge of general high-frequency words.

### **What is the value of the ASWL for L2 vocabulary research, learning, and teaching?**

The ASWL is a valuable resource for L2 vocabulary research, learning, and teaching. For research, the ASWL makes a number of contributions. First, the approach taken to develop the list is innovative. It makes the best use of the two approaches toward identifying academic words and helps to deal with the challenge of distinguishing academic and non-academic vocabulary to some extent. Following Gardner and Davies's (2014) approach, the ASWL was created from scratch, which avoids the limitations related to ready-made lists of general high-frequency vocabulary and allows for differences between academic speech and general conversation to some extent (Csomay 2006). Creating the ASWL from scratch means not including 714 general

high-frequency words that did not meet the selection criteria (e.g., delicious, pudding, grin) and including 455 words at lower frequency levels that satisfied the criteria (e.g., theory, define, factor). This results in a list with a more attainable size and higher coverage than a general high-frequency wordlist. Learners' attention will be directed to the most important words in academic speech, especially the items whose occurrences in general conversation are not frequent enough for incidental learning to happen. However, Gardner and Davies (2014) only included general high-frequency words in their AVL if the frequency of these words was at least 50% higher in academic texts than in non-academic texts (ratio of 1.5). In contrast, the ASWL includes general high-frequency words if they have high frequency and wide range in academic texts irrespective of the ratio between their frequency in academic texts and in non-academic texts. Thus, the present study gives credit to items which have high frequency and wide range in both academic and non-academic texts, because they are also important for learners' comprehension of academic speech. Examples of these items are investigate, propose, think, and issue, which appeared in all 24 subject areas and had a dispersion from 0.75 to 0.91 in the first academic spoken corpus. The ratio between their frequency in the first academic spoken corpus and non-academic spoken corpus was 0.56 (investigate), 0.60 (propose), 0.61 (think), and 1.4 (issue).

Following Xue and Nation (1984), Coxhead (2000), and Nesi (2002), the development of the ASWL considered learners' existing knowledge of general high-frequency vocabulary. However, instead of making a benchmark of the number of general words that all learners should acquire before learning the ASWL, the ASWL was divided into four levels according to the word frequency levels so that the list can be well-suited to a wide range of learners. No previous research on academic wordlists has looked at the issue from this angle. Therefore, this study effectively updates Nation and Hwang's (1995) question about the point at which learners should study specialized vocabulary. It suggests that there is no clear-cut boundary between general service vocabulary and specialized vocabulary, and learners can start learning specialized vocabulary at any stage to match their existing vocabulary knowledge. Moreover, focusing on the lexical items that are beyond learners' existing vocabulary levels and creating opportunities for learners to repeatedly meet and use these words in texts from their specific subject areas can help learners to be better aware of the discipline-specific meanings of these words. Additionally, grading the ASWL into levels makes the list more adaptable to learners' proficiencies and allows teachers and learners to avoid repeatedly learning known items. In brief, the development of the

ASWL provides subsequent studies into specialized vocabulary with a useful way to make corpus-based wordlists more suitable to L2 learning and teaching.

Second, due to the lack of academic spoken wordlists, previous research used Coxhead's (2000) AWL to represent academic words in academic speech (e.g., Vidal, 2003, 2011). Yet, the nature of academic spoken vocabulary is different from that in written discourse (Dang & Webb, 2014; Biber, 2006). Therefore, the development of the ASWL provides a useful basis for future research on academic spoken vocabulary in multiple aspects such as vocabulary load, incidental learning, testing, and technical wordlist development. Moreover, the present study provides clearer descriptions of the corpus development and word selection criteria than any previous research on developing specialized wordlists. This then allows other researchers to replicate this study more easily.

As for those involving in L2 vocabulary learning and teaching, the ASWL helps learners, teachers, course designers, and material writers to set vocabulary learning goals and learning sequences to enhance learners' comprehension of academic speech. For example, at the beginning of an English language program, learners' vocabulary levels can be measured by Webb, Sasao, and Ballance's' (n.d.) New Vocabulary Levels Test. Depending on learners' existing levels of general vocabulary and their specific learning and teaching contexts, learners and teachers can use the sequence in Figure 1 as a guide to set long-term learning goals for the learners. The BNC/COCA lists were used to guide the learning sequence because the 1<sup>st</sup> and 2<sup>nd</sup> 1,000 BNC/COCA word lists are the most suitable general high-frequency word lists for L2 learners (Dang, 2017; Dang & Webb, 2016a). Moreover, the BNC/COCA lists have been widely used to represent general vocabulary in numerous vocabulary studies. Using these lists to guide the learning sequence, therefore, makes it possible for teachers, course designers, and material writers to incorporate the findings of the present study with the findings of other studies related to the BNC/COCA lists to organize a more effective vocabulary learning program. The sequence shown in Figure 1 also helps learners and teachers avoid repeatedly teaching and learning known items, and makes it easier for teachers and course designers to incorporate teaching the ASWL to groups of learners with different vocabulary levels and learning purposes. Let us take learners who have not mastered the most frequent 1,000 BNC/COCA word-families as an example. If they want to go straight to the lexical items that appear frequently in their academic studies, they

can start learning the ASWL at Level 1. However, if they would like to acquire the items that are useful for general conversation before moving to the words occurring frequently in their academic disciplines, they can learn items from the BNC/COCA wordlists first. Once they are satisfied with their levels of general vocabulary, they can start learning the ASWL at the level which is relevant to their general vocabulary level at that time.

[FIGURE 1 NEAR HERE]

While the ASWL levels provide guidance in setting long-term learning goals for L2 learners, its frequency ranked sub-lists of manageable size make it easier for learners, teachers, and course designers to set short-term learning goals and incorporate the ASWL in their language learning program (Dang & Webb, 2016b; Nation, 2016). Moreover, teaching and learning the ASWL according to the rank order of sub-lists ensures that the most useful items are learned first as well as allowing programs to prepare a curriculum that covers all sub-lists, and avoids teaching the same items between courses (Coxhead, 2000).

Distinguishing between function words and lexical words takes into account the difference in the way these words are learned (Dang & Webb, 2016b; Nation, 2016). As lexical words are more salient than function words in a text, the way to deal with lexical words should be different from the way to deal with most function words (Carter & McCarthy, 1988). It will be best to sequence the teaching of lexical words according to their frequency. Yet, it is more reasonable to incorporate teaching function words with other components of language lessons because of their lack of salience in the text. No previous academic wordlists make a distinction between function words and lexical words and are adaptable to learners' proficiencies, which makes the ASWL more pedagogically appropriate.

A few issues should be noted in the implementation of the ASWL. First, the ASWL levels and sub-lists should be treated as a guide rather than a handbook for learners and teachers to strictly follow. Second, the ASWL consists of 1,741 word-families. The rationale for using word-families is that the learning burden of a word (e.g., sadly) which is morphologically related to a known word (e.g., sad) may be less than that of an unrelated word (e.g., abolish). This assumption has strong evidence from psycholinguistic studies with L1 children (Bertram, Baayen, & Schreuder, 2000; Nagy, Anderson, Schommer, Scott, & Stallman, 1989) and L2 learners (Mochizuki & Aizawa, 2000; Schmitt & Meara, 1997; Sasao & Webb, 2017; Schmitt &

Zimmerman, 2002). However, as learners' knowledge of affixes (or derivational knowledge) increases incrementally, it is important for learners to keep in mind that they should not restrict their learning to 1,741 ASWL headwords only but should expand their knowledge of their family members and knowledge of affixes. Similarly, the ASWL and BNC/COCA2000 are presented in the format of wordlists, but this does not mean that these lists should be learned and taught solely by decontextualized methods (Coxhead, 2000). Knowing a word involves many types of vocabulary knowledge (Nation, 2013). Therefore, once the learning goals have been set, it is important for teachers and material designers to design learning activities and materials for learners to repeatedly encounter and use these words in different contexts related to their target subject areas (Coxhead, 2000). In this way, learners can acquire, consolidate, and expand on their knowledge of these words in a meaningful way. Nation's (2007) four strands provides a useful framework for organizing learning opportunities. Third, although 90% and 95% coverage of academic spoken English are important goals for EAP learners, once they have achieved these goals, they should continue to expand their vocabulary knowledge to reach higher coverage figures such as 98%. It is because the higher coverage, the better comprehension (Schmitt, Jiang, & Grabe, 2011). Finally, the present study acknowledges that to achieve academic success, learners need to have a good knowledge of both academic spoken and written vocabulary. That said, many resources are available for students to improve their knowledge of academic written vocabulary, but academic spoken vocabulary resources are very limited.

## **Limitations**

The present study adopted several assumptions from earlier studies that could be questioned. The coverage figures in this study were based on the analysis of written transcription of spoken data with the assumption that learners are able to recognize the spoken forms of the words as well as proper nouns and marginal words. This assumption has been adopted by previous corpus-driven research investigating the vocabulary load of spoken texts (e.g., Nation, 2006). However, it should be noted that, although L2 learners' aural and orthographic knowledge are closely related (Milton, 2009), the gap between the two kinds of knowledge may vary. Kobeleva (2012) found that previous knowledge of proper names has significant importance in listening comprehension; therefore, it may be optimistic to assume that proper nouns do not require previous knowledge when being encountered in listening (Nation, 2016). A second assumption is the use of Juillard

and Chang-Rodrigues's (1964) D for dispersion. Biber, Reppen, Schnur, and Ghanem (2016) found a large decrease in the sensitivity of D when the computations were based on a large number of corpus parts. As a result, there exists the possibility that the ASWL contains some unevenly distributed items despite meeting the Juilland's D criterion. Like most previous research on developing academic wordlists (e.g., Gardner & Davies, 2014; Xue & Nation, 1984), this study only looked at the influence of lexical coverage on comprehension. However, there are many other factors affecting listening comprehension such as background knowledge (Schmidt-Rinehart, 1994), speech rate (Flowerdew & Miller, 1992), and interaction (Flowerdew & Miller, 1992). This study did not look at the extent to which each member occurs but rather the extent to which the unit of counting occurs, which means chances that a member of a word-family might not occur. Finally, the second academic spoken corpus and the academic written corpus are not as well-balanced as the first academic spoken corpus used to develop the ASWL.

## **Future research**

One direction for future research is further validation and application of the ASWL. This study was based on hypothetical calculation that the ASWL can assist learners with different proficiency levels to reach a larger coverage of academic spoken English. Intervention research with real learners can provide further insight into the actual coverage learners may gain from the ASWL. Given the importance of multi-words in academic speech (Simpson-Vlach & Ellis, 2010), further research should look at which words commonly collocate with the ASWL words. Moreover, Webb et al.'s (n.d.) NVLT was used to identify the relevant ASWL levels for learners to focus on. Unlike Schmitt, Schmitt, and Clapham's (2001) Vocabulary Levels Test (VLT), the NVLT was based on the BNC/COCA lists and has separate 1<sup>st</sup> and 2<sup>nd</sup> 1,000 word levels. Therefore, it can provide more precise information about learners' vocabulary knowledge at each 1,000 word frequency level. However, like Schmitt et al.'s (2001) VLT, the NVLT is not an aural test. There is a need for further development of tests designed to measure both knowledge of the ASWL, as well as tests designed to measure the spoken forms of the ASWL and BNC/COCA words, and these tests should be based on rigorous validations.

Further research on developing spoken wordlists for each disciplinary group would be valuable (Hyland & Tse, 2007). These lists may better serve the needs of ESP or ESAP classes with learners planning to study the same or similar disciplines. Another option is to develop an



academic written word list (AWWL) that follows the same approach as the ASWL, to allow a valid comparison between the most frequent and wide ranging items in academic spoken and written discourse. Given that L2 learners need to comprehend both academic spoken and written texts for their success, the AWWL and ASWL together might provide L2 learners with better support to succeed in both spoken and written interactions. The development of these specialized wordlists might make use of a powerful measure of dispersion other than Juilland and Chang-Rodrigues's (1964) D (e.g., Gries's (2008) DP), and use validating corpora mirroring the corpora from which these lists are developed. It would also be useful to investigate the extent to which using different dispersion criteria leads to different lists, and the degree to which a list produced from the validating academic spoken corpus is similar to that in the source academic spoken corpus.

Future research could examine other features of vocabulary in academic spoken English. This study considers lectures, seminars, labs, and tutorials as a whole. As different speech events have distinctive linguistic features (Biber, 2006), it is useful for future research to examine the vocabulary in each of these speech events as well as in other kinds of speech events (e.g., office hours, conference presentations).

## **Conclusion**

The ASWL is an example of how corpus-based wordlists can better support the continual vocabulary development of L2 learners irrespective of their proficiencies, disciplines, and institutional structures. It contains 1,741 word-families with high frequency, wide range, and even distribution in academic speech. Truly reflecting the most frequent and wide ranging vocabulary in academic spoken English, the ASWL may help learners to reach 92%-96% coverage of academic speech. The ASWL levels and sub-lists are useful resources for setting learning goals and sequences as well as designing courses and materials to enhance L2 learners' comprehension of academic spoken English. The method used to develop the ASWL provides a useful direction for future research on specialized wordlists. Importantly, the ASWL provides a basic foundation for further research into academic spoken vocabulary.

## **Acknowledgements**

We would like to thank the reviewers and editor for their useful feedback, and to the following publishers and researchers for their generosity in letting us use their materials to create our corpora: Cambridge University Press, Pearson, Lynn Grant, the lecturers at Victoria University of Wellington, the researchers in the British Academic Spoken English corpus project, the British Academic Written English corpus project, the International Corpus of English project, the Massachusetts Institute of Technology Open courseware project, the Open American National corpus project, the Santa Barbara Corpus of Spoken American-English project, the Stanford Engineering Open courseware project, the University of California, Berkeley Open courseware project, and the Yale University Open courseware project.

## **Author note**

### **Thi Ngoc Yen Dang (corresponding author)**

Vietnam National University  
84 PhamThe Hien, TO 12, Tran Hung Dao  
Thai Binh City, Thai Binh, Vietnam  
[ngocyen1011@gmail.com](mailto:ngocyen1011@gmail.com)

### **Averil Coxhead**

School of Linguistics and Applied Language Studies  
Victoria University of Wellington  
PO Box 600, Wellington, New Zealand  
[averil.coxhead@vuw.ac.nz](mailto:averil.coxhead@vuw.ac.nz)

### **Stuart Webb**

Faculty of Education  
University of Western Ontario  
London, Ontario, Canada  
[swebb27@uwo.ca](mailto:swebb27@uwo.ca)

## References

- Adolphs, S., & Carter, R. (2013). *Spoken corpus linguistics: From monomodal to multimodal*. New York: Routledge.
- Aguilar, M. (2016). Seminars. In K. Hyland & P. Shaw (Eds.), *The Routledge handbook of English for Academic Purposes* (pp. 335–347). London: Routledge.
- Banister, C. (2016). The Academic Word List: Exploring teacher practices, attitudes and beliefs through a web-based survey and interviews. *The Journal of Teaching English for Specific and Academic Purposes*, 4(2), 309–325.
- Bauer, L., & Nation, P. (1993). Word families. *International Journal of Lexicography*, 6(4), 253–279. doi:10.1093/ijl/6.4.253
- Becher, T. (1989). *Academic tribes and territories*. Bristol: The Society for Research into Higher Education and Open University Press.
- Becker, A. (2016). L2 students' performance on listening comprehension items targeting local and global information. *Journal of English for Academic Purposes*, 24, 1–13. doi:10.1016/j.jeap.2016.07.004
- Berman, R., & Cheng, L. (2001). English academic language skills: Perceived difficulties by undergraduate and graduate students, and their academic achievement. *Canadian Journal of Applied Linguistics*, 4(1–2), 25–40.
- Bertram, R., Baayen, R. H., & Schreuder, R. (2000). Effects of family size for complex words. *Journal of Memory and Language*, 43(3), 390–405. <https://doi.org/10.1006/jmla.1999.2681>
- Biber, D. (1993). Representativeness in corpus design. *Literary and Linguistic Computing*, 8(4), 243–257. doi:10.1093/lc/8.4.243
- Biber, D. (2006). *University language: A corpus-based study of spoken and written registers*. Amsterdam: John Benjamins Publishing.

- Biber, D., Conrad, S., Reppen, R., Byrd, P., & Helt, M. (2002). Speaking and writing in the university: A multidimensional comparison. *TESOL Quarterly*, 36(1), 9–48. doi: [10.2307/3588359](https://doi.org/10.2307/3588359)
- Biber, D., Johansson, S., Leech, G., Conrad, S., & Finegan, E. (1999). *Longman grammar of spoken and written English*. London: Longman.
- Biber, D., Reppen, R., Schnur, E., & Ghanem, R. (2016). On the (non)utility of Juilland’s D to measure lexical dispersion in large corpora. *International Journal of Corpus Linguistics*, 21(4), 439–464. doi: [10.1075/ijcl.21.4.01bib](https://doi.org/10.1075/ijcl.21.4.01bib)
- Brezina, V., & Gablasova, D. (2015). Is there a core general vocabulary? Introducing the New General Service List. *Applied Linguistics*, 36(1), 1–22. doi: [10.1093/applin/amt018](https://doi.org/10.1093/applin/amt018)
- Carter, R., & McCarthy, M. (1988). *Vocabulary and language teaching*. London: Longman.
- Coxhead, A. (2000). A new academic word list. *TESOL Quarterly*, 34(2), 213–238. doi: [10.2307/3587951](https://doi.org/10.2307/3587951)
- Coxhead, A. (2011). The Academic Word List 10 years on: Research and teaching implications. *TESOL Quarterly*, 45(2), 355–362. doi: [10.5054/tq.2011.254528](https://doi.org/10.5054/tq.2011.254528)
- Coxhead, A. (2016). Reflecting on Coxhead (2000), “A New Academic Word List.” *TESOL Quarterly*, 50(1), 181–185. doi: [10.1002/tesq.287](https://doi.org/10.1002/tesq.287)
- Coxhead, A., & Hirsh, D. (2007). A pilot science-specific word list. *Revue Française de Linguistique Appliquée*, 12(2), 65–78.
- Chung, T. M., & Nation, P. (2004). Identifying technical vocabulary. *System*, 32(2), 251–263. <https://doi.org/10.1016/j.system.2003.11.008>
- Crossley, S., Salsbury, T., & McNamara, D. (2010). The development of polysemy and frequency use in English Second Language speakers. *Language Learning*, 60(3), 573–605. doi: [10.1111/j.1467-9922.2010.00568.x](https://doi.org/10.1111/j.1467-9922.2010.00568.x)
- Csomay, E. (2006). Academic talk in American university classrooms: Crossing the boundaries

- of oral-literate discourse? *Journal of English for Academic Purposes*, 5(2), 117–135. doi:10.1016/j.jeap.2006.02.001
- Dang, T. N. Y. (2017). Investigating vocabulary in academic spoken English: Corpora, teachers, and learners (Unpublished PhD thesis). Victoria University of Wellington, Wellington, New Zealand.
- Dang, T. N. Y., & Webb, S. (2014). The lexical profile of academic spoken English. *English for Specific Purposes*, 33, 66–76. doi: 10.1016/j.esp.2013.08.001
- Dang, T. N. Y., & Webb, S. (2016a). Evaluating lists of high-frequency words. *ITL – International Journal of Applied Linguistics*, 167(2), 132–158. doi: 10.1075/itl.167.2.02dan
- Dang, T. N. Y., & Webb, S. (2016b). Making an essential word list. In I. S. P. Nation (Ed.), *Making and using word lists for language learning and testing* (pp. 153–167). Amsterdam: John Benjamins.
- Flowerdew, J., & Miller, L. (1992). Student perceptions, problems and strategies in second language lecture comprehension. *RELC*, 23(2), 60–80. doi:10.1177/003368829202300205
- Francis, W. N., & Kučera, H. (1982). *Frequency analysis of English usage: Lexicon and grammar*. Boston: Houghton Mifflin.
- Gardner, D., & Davies, M. (2014). A new academic vocabulary list. *Applied Linguistics*, 35(3), 305–327. doi:10.1093/applin/amt015
- Gries, S. T. (2008). Dispersions and adjusted frequencies in corpora. *International Journal of Corpus Linguistics*, 13(4), 403–437. <https://doi.org/10.1075/ijcl.13.4.02gri>
- Henriksen, B., & Danelund, L. (2015). Studies of Danish L2 learners' vocabulary knowledge and the lexical richness of their written production in English. In P. Pietilä, K. Doró, & R. Pipalová(Eds.), *Lexical issues in L2 writing* (pp. 1–27). Newcastle upon Tyne: Cambridge ScholarsPublishing.
- Hyland, K. (2000). *Disciplinary discourses: Social interactions in academic writing*. London: Longman.

- Hyland, K. (2016). General and specific EAP. In K. Hyland & P. Shaw (Eds.), *The Routledge handbook of English for Academic Purposes* (pp. 17–29). London: Routledge.
- Hyland, K., & Tse, P. (2007). Is there an “academic vocabulary”? *TESOL Quarterly*, 41(2), 235–253. doi:10.1002/j.1545-7249.2007.tb00058.x
- Jones, W. A. (2011). Variation among academic disciplines: An update on analytical frameworks and research. *The Journal of the Professoriate*, 6(1), 9–27.
- Juilland, A. G., & Chang-Rodríguez, E. (1964). *Frequency dictionary of Spanish words*. London: Mouton.
- Kobeleva, P. P. (2012). Second language listening and unfamiliar proper names: Comprehension barrier? *RELC Journal*, 43(1), 83–98. doi:10.1177/0033688212440637
- Laufer, B. (1998). The development of passive and active vocabulary in a second language: Same or different? *Applied Linguistics*, 19(2), 255–271. doi:10.1093/applin/19.2.255
- Leech, G. N., Rayson, P., & Wilson, A. (2001). *Word frequencies in written and spoken English*. Harlow: Longman.
- Lei, L., & Liu, D. (2016). A new medical academic word list: A corpus-based study with enhanced methodology. *Journal of English for Academic Purposes*, 22, 42–53. doi:10.1016/j.jeap.2016.01.008
- Liu, J., & Han, L. (2015). A corpus-based environmental academic word list building and its validity test. *English for Specific Purposes*, 39, 1–11. doi:10.1016/j.esp.2015.03.001
- Lynch, T. (2011). Academic listening in the 21st century: Reviewing a decade of research. *Journal of English for Academic Purposes*, 10(2), 79–88. doi:10.1016/j.jeap.2011.03.001
- MacDonald, M., Badger, R., & White, G. (2000). The real thing?: Authenticity and academic listening. *English for Specific Purposes*, 19(3), 253–267. doi:10.1016/S0889-4906(98)00028-

- Matthews, J., & Cheng, J. (2015). Recognition of high frequency words from speech as a predictor of L2 listening comprehension. *System*, 52, 1–13. doi:10.1016/j.system.2015.04.015
- Mauranen, A. (2004). Speech corpora in the classroom. In G. Aston, S. Bernardini, & D. Stewart (Eds.), *Corpora and language learners* (pp. 197–213). Amsterdam: John Benjamins.
- McCarthy, M., & Carter, R. (1997). Written and spoken vocabulary. In N. Schmitt & M. McCarthy (Eds.), *Vocabulary: Description, acquisition and pedagogy* (pp. 20–39). Cambridge: Cambridge University Press.
- Milton, J. (2009). *Measuring second language vocabulary acquisition*. Bristol: Multilingual Matters.
- Mochizuki, M., & Aizawa, K. (2000). An affix acquisition order for EFL learners: An exploratory study. *System*, 28(2), 291–304. doi:10.1016/S0346-251X(00)00013-0
- Mulligan, D., & Kirkpatrick, A. (2000). How much do they understand? Lectures, students and comprehension. *Higher Education Research & Development*, 19(3), 311–335. doi:10.1080/758484352
- Nagy, W., Anderson, R. C., Schommer, M., Scott, J. A., & Stallman, A. C. (1989). Morphological families in the internal lexicon. *Reading Research Quarterly*, 24(3), 262–282. doi: 10.2307/747770
- Nation, I. S. P. (2006). How large a vocabulary is needed for reading and listening? *Canadian Modern Language Review*, 63(1), 59–82. doi: 10.3138/cmlr.63.1.59
- Nation, I. S. P. (2007). The four strands. *Innovation in Language Learning and Teaching*, 1(1), 1–12. doi: 10.2167/illt039.0
- Nation, I. S. P. (2012). The BNC/COCA word family lists. Retrieved from <http://www.victoria.ac.nz/lals/about/staff/paul-nation>

- Nation, I. S. P. (2013). *Learning vocabulary in another language* (2nd ed.). Cambridge: Cambridge University Press.
- Nation, I. S. P. (2016). *Making and using word lists for language learning and testing*. Amsterdam: John Benjamins.
- Nation, I. S. P., Heatley, A., & Coxhead, A. (2002). Range: A program for the analysis of vocabulary in texts. Retrieved from <http://www.vuw.ac.nz/lals/staff/paul-nation/nation.aspx>
- Nation, P., & Hwang, K. (1995). Where would general service vocabulary stop and special purposes vocabulary begin? *System*, 23(1), 35–41. doi: 10.1016/0346-251X(94)00050-G
- Nation, I. S. P., & Webb, S. (2011). *Researching and analyzing vocabulary*. Boston: Heinle, Cengage Learning.
- Nation, P. (2004). A study of the most frequent word families in the British National Corpus. In P. Bogaards & B. Laufer (Eds.), *Vocabulary in a second language: Selection, acquisition, and testing* (pp. 3–13). Amsterdam: John Benjamins.
- Nesi, H. (2002). An English Spoken Academic Word List. In A. Braasch & C. Povlsen (Eds.), *Proceedings of the Tenth EURALEX International Congress* (Vol. 1, pp. 351–358). Copenhagen, Denmark. Retrieved from [http://www.euralex.org/elx\\_proceedings/Euralex2002/036\\_2002\\_V1\\_Hilary%20Nesi\\_An%20English%20Spoken%20Academic%20Wordlist.pdf](http://www.euralex.org/elx_proceedings/Euralex2002/036_2002_V1_Hilary%20Nesi_An%20English%20Spoken%20Academic%20Wordlist.pdf)
- Neumann, R. (2001). Disciplinary differences and university teaching. *Studies in Higher Education*, 26(2), 135–146. doi: 10.1080/03075070120052071
- Nguyen, T. M. H., & Webb, S. (2016). Examining second language receptive knowledge of collocation and factors that affect learning. *Language Teaching Research*, 1 –23. <https://doi.org/10.1177/1362168816639619>
- Nurweni, A., & Read, J. (1999). The English vocabulary knowledge of Indonesian university



- students. *English for Specific Purposes*, 18(2), 161–175. doi:10.1016/S0889-4906(98)00005-2
- O’Keeffe, A., McCarthy, M., & Carter, R. (2007). *From corpus to classroom: Language use and language teaching*. Cambridge: Cambridge University Press.
- Pinchbeck, G. G. (2014). Lexical frequency profiling of a large sample of Canadian high school diploma exam expository writing: L1 and L2 academic English. Presented at the Roundtable presentation at American Association of Applied Linguistics, Portland, OR, USA.
- Sasao, Y., & Webb, S. (2017). The Word Part Levels Test. *Language Teaching Research*, 21(1), 12–30. <https://doi.org/10.1177/1362168815586083>
- Schmidt-Rinehart, B. C. (1994). The effects of topic familiarity on second language listening comprehension. *The Modern Language Journal*, 78(2), 179–189. doi:10.1111/j.1540-4781.1994.tb02030.x
- Schmitt, N., Cobb, T., Horst, M., & Schmitt, D. (2015). How much vocabulary is needed to use English? Replication of van Zeeland & Schmitt (2012), Nation (2006) and Cobb (2007). *Language Teaching*. doi:10.1017/S0261444815000075
- Schmitt, N., Jiang, X., & Grabe, W. (2011). The percentage of words known in a text and reading comprehension. *The Modern Language Journal*, 95(1), 26–43. doi:10.1111/j.1540-4781.2011.01146.x
- Schmitt, N., & Meara, P. (1997). Research vocabulary through a word knowledge framework. *Studies in Second Language Acquisition*, 19(1), 17–36. doi:10.1017/S0272263197001022
- Schmitt, N., & Schmitt, D. (2014). A reassessment of frequency and vocabulary size in L2 vocabulary teaching. *Language Teaching*, 47(4), 484–503. doi: 10.1017/S0261444812000018
- Schmitt, N., Schmitt, D., & Clapham, C. (2001). Developing and exploring the behaviour of two new versions of the Vocabulary Levels Test. *Language Testing*, 18(1), 55–88. doi:10.1191/026553201668475857

- Schmitt, N., & Zimmerman, C. B. (2002). Derivative word forms: What do learners know? *TESOL Quarterly*, 36(2), 145–171. doi:10.2307/3588328
- Shin, D., & Nation, P. (2008). Beyond single words: the most frequent collocations in spoken English. *ELT Journal*, 62(4), 339–348. doi: 10.1093/elt/ccm091
- Simpson-Vlach, R., & Ellis, N. C. (2010). An Academic Formulas List: New methods in phraseology research. *Applied Linguistics*, 31(4), 487–512. doi:10.1093/applin/amp058
- Sinclair, J. (1991). *Corpus, concordance, collocation*. Oxford: Oxford University Press.
- Thompson, P. (2006). A corpus perspective on the lexis of lectures, with a focus on economics lectures. In K. Hyland & M. Bondi (Eds.), *Academic discourse across disciplines* (pp. 253–270). New York: Peter Lang.
- Townsend, D., & Collins, P. (2009). Academic vocabulary and middle school English learners: An intervention study, 22(9), 993–1019. doi:10.1007/s11145-008-9141-y
- van Zeeland, H., & Schmitt, N. (2013). Lexical coverage in L1 and L2 listening comprehension: The same or different from reading comprehension? *Applied Linguistics*, 34(4), 457–479. doi:10.1093/applin/ams074
- Vidal, K. (2003). Academic listening: A source of vocabulary acquisition? *Applied Linguistics*, 24(1), 56–89. doi:10.1093/applin/24.1.56
- Vidal, K. (2011). A comparison of the effects of reading and listening on incidental vocabulary acquisition. *Language Learning*, 61(1), 219–258. doi:10.1111/j.1467-9922.2010.00593.x
- Wang, J., Liang, S., & Ge, G. (2008). Establishment of a Medical Academic Word List. *English for Specific Purposes*, 27(4), 442–458. doi:10.1016/j.esp.2008.05.003
- Ward, J., & Chuenjundaeng, J. (2009). Suffix knowledge: Acquisition and applications. *System*, 37(3), 461–469. doi: 10.1016/j.system.2009.01.004

Webb, S. A., & Chang, A. C.-S. (2012). Second language vocabulary growth. *RELC Journal*, 43(1), 113–126. doi: 10.1177/0033688212439367

Webb, S., Sasao, Y., & Ballance, O. (n.d.). New Vocabulary Levels Test. Retrieved from at [http://vuw.qualtrics.com/jfe/form/SV\\_4MG1wByukg1JoTb](http://vuw.qualtrics.com/jfe/form/SV_4MG1wByukg1JoTb)

West, M. (1953). *A general service list of English words*. London: Longman, Green.

Wei, Z. (2014). Does teaching mnemonics for vocabulary learning make difference? Putting the keyword method and the word part technique to the test. *Language Teaching Research*, 19(1), 43–69. <https://doi.org/10.1177/1362168814541734>

Xue, G., & Nation, I. S. P. (1984). A university word list. *Language Learning and Communication*, 3(2), 215–229.

Yang, M.-N. (2015). A nursing academic word list. *English for Specific Purposes*, 37, 27–38. doi: 10.1016/j.esp.2014.05.003

Zipf, G. K. (1935). *The psycho-biology of language*. Cambridge: MIT Press.

**Table 1** Academic and non-academic corpora for the ASWL study

Corpus	Purpose	Size (running-words)
1 <sup>st</sup> academic spoken corpus	Develop the ASWL	13,029,661
2 <sup>nd</sup> academic spoken corpus	Determine if the ASWL accurately reflects the vocabulary in academic speech	12,740,619
Academic written corpus	Examine if the ASWL reflects spoken vocabulary	13,449,584
Non-academic spoken corpus	Examine if the ASWL reflects academic vocabulary	13,863,628

**Table 2** Bauer and Nation’s (1993) word family levels (adapted from Nation, 2016, p.27)

Word family level	Affixes
-------------------	---------

Level 1	A different form is a different word
Level 2	Inflectional suffixes: plural, third person singular present tense, past tense, past participle, -ing, comparative, superlative, and possessive (8 affixes)
Level 3	Most frequent and regular derivational affixes: -able,-er,-ish,-less,-ly,-ness,-th,-y, non-, un- (all with restricted use) (10 affixes)
Level 4	Frequent, orthographically regular affixes: -al,-ation, -ess, -ful, -ism, -ist, -ity, -ize,-ment, -ous, in- (all with restricted use) (11 affixes)
Level 5	Regular but infrequent affixes: -age, -al, -ally, -an,-ance, -ant, -ary, -atory, -dom, -eer, -en, -en, -ence, -ent, ery, -ese, -esque, -ette, -hood, -i, -ian, -ite, -let, -ling, -ly, -most, -ory,-ship, -ward, -ways, -wise, ante-, anti-, arch-, bi-, circum-, counter, en-, ex-, fore-, hyper-, inter-, mid-, mis-, neo-, post-, pro-, semi-, sub-, un- (50 affixes)
Level 6	Frequent but irregular affixes: -able, -ee, -ic, -ify, -ion, -ist, -ition, -ive (ative), -th, -y, pre-, re- (12 affixes)
Level 7	Classical roots and affixes

**Table 3** Lexical profile of the ASWL

ASWL level	BNC/COCA word level	Number of word-families	Coverage (%)	Examples
Level 1	1 <sup>st</sup> 1,000	830	81.62	alright, know, stuff
Level 2	2 <sup>nd</sup> 1,000	456	5.23	therefore, determine, approach
Level 3	3 <sup>rd</sup> 1,000	380	2.85	achieve, significant, aspect, review
Level 4	4 <sup>th</sup> 1,000	49	0.28	straightforward, triangle, differentiate
	5 <sup>th</sup> 1,000	13	0.08	arbitrary, coefficient, analytic
	6 <sup>th</sup> 1,000	6	0.03	radius, optimise, intuition
	7 <sup>th</sup> 1,000	3	0.01	subtract, gamma, inverse
	8 <sup>th</sup> 1,000	1	0.004	theorem

9 <sup>th</sup> 1,000	1	0.01	exponential
10 <sup>th</sup> 1,000	1	0.01	semester
Outside BNC/COCA	1	0.01	so-called
<b>Total</b>	<b>1,741</b>	<b>90.13</b>	

**Table 4** Coverage of Coxhead's (2000) AWL, Gardner and Davies's (2014) AVL and the ASWL in the two academic spoken corpora

Word lists	Number of Level-6 word-families	Coverage (%)	
		1 <sup>st</sup> academic spoken corpus	2 <sup>nd</sup> academic spoken corpus
AWL	570	4.17	4.03
AVL	1,983	23.88	23.78
ASWL	1,741	90.13	89.59

**Table 5** Potential coverage gained by learners with the aid of the ASWL (%)

Existing vocabulary level (BNC/COCA word- families)	Number of ASWL word-families beyond learners' level	Without proper nouns & marginal words		With proper nouns & marginal words	
		1st academic spoken corpus	2 <sup>nd</sup> academic spoken corpus	1st academic spoken corpus	2nd academic spoken corpus
Less than 1,000	1,741	90.13	89.59	92.60	92.35
1,000	911	90.79	90.12	93.26	92.88
2,000	455	91.89	91.03	94.36	93.79
3,000	75	93.33	92.24	95.80	95.0
Proper nouns		1.23	1.40		
Marginal words		1.24	1.36		