## Book Section:

**Dang, T. N. Y**., & Webb, S. (2016). Making an essential word list for beginners. In I. S. P. Nation (Ed.), *Making and using word lists for language learning and testing* (pp. 153–167). Amsterdam: John Benjamins.

# Making an essential word list for beginners

*Thi Ngoc Yen Dang and Stuart Webb*

This chapter describes a word list study which expands on earlier studies and creates a practical wordlist that would provide a starting point for L2 beginners' lexical development. An initial aim is to identify which words should be included in an essential wordlist for L2 beginners. A second aim is to determine how many items should be included in a wordlist for L2 beginners using three criteria: practicability, change in the coverage curve, and amount of lexical coverage. The word list could serve as the foundation for L2 beginner lexical development. The points to note about the study are its choice of the unit of counting, the size of the list and its sub-lists, the treatment of proper noun homonyms and the extensive validation of the list.

**Which items should be included in a word list for beginners?**

Analysis of established lists that were developed from large corpora using precise and valid methodologies may provide a reliable list of essential vocabulary for beginners. West's (1953) GSL was chosen as one of the source lists in the present study because it is the oldest and most influential high-frequency wordlist. Nation's (2006) BNC2000, Nation's (2012) BNC/COCA2000, and Brezina and Gablasova's (2013) New-GSL were chosen because they are three recently-created high-frequency wordlists, and earlier studies (Dang & Webb, under review; Brezina & Gablasova, 2013) have shown that these lists provided higher lexical coverage than the GSL in multiple corpora. Another high-frequency wordlist, Browne's (2013) New General Service List, was created recently. It was not used in the present study for two reasons. First, preliminary analysis of the list as a whole showed that Browne's (2013) list provided lower average coverage per item than any of the four lists in the present study including the GSL. Second, there has been little written about the way it was developed.

In a comparison of the lexical coverage provided by items in the GSL, BNC2000, BNC/COCA2000, and New-GSL in nine spoken and nine written corpora, Dang and Webb (under review) found that each list had both strong and weak items. This suggests that a list made of the best items from the four lists may provide greater coverage than any one list. Moreover, because high-frequency words occur frequently in a wide range of texts, validation of the items in high-frequency wordlists should be based on coverage in a larger number of corpora with a greater degree of variation of English language than has been used in the earlier studies. The present study aims to fill this gap by ranking the items from all lists based on their lexical coverage in a large number of corpora representing different discourse types and varieties of English.

**What should be the unit of counting in a wordlist created for beginners?**

An important issue when developing wordlists is the unit of counting. The GSL, BNC2000, and BNC/COCA2000 used Level 6 word-families (Bauer & Nation, 1993) as the unit of counting while the New-GSL used lemmas (Level 2 word families). The choice of word-families as the unit of counting is based on the assumption that, if learners know one word-form, they may recognize its inflected and closely derived forms. In contrast to Level 6 word families, the choice of Level 2 families is based on the

assumption that, if learners know one word-form, they may only recognize its inflected forms. Level 2 word families in the Bauer and Nation scheme involve only inflected forms. Each option has its advantages and disadvantages, and the choice of unit of counting should be based on the characteristics of target list-users (Gardner, 2007).

In wordlists for L2 beginners, Level 2 word-families are more suitable than Level 6 word-families for two reasons. First, L2 beginners' morphological awareness may be limited, and it may be inappropriate to assume that if they know one member of a word-family, they may recognize its derivational forms. This is supported by Schmitt and Zimmerman's (2002) and Ward and Chuenjundaeng's (2009) studies that found that not all derivational members of a word-family were known by L2 learners. Second, for L2 beginners who lack sufficient English morphological knowledge and their teachers, a Level 2 word family list might be more useful than a Level 6 word-family list. Level 2 lists consist of mainly high-frequency lemmas (*study*) while word-family lists are made up of both high-frequency (*study*) and low-frequency lemmas (*studious, studiously*). Introducing Level 2 lists to L2 beginners will draw their attention to the high-frequency words first. By developing knowledge of these most important forms, it may be easier to learn the infrequent members from the same word-family at a later stage of lexical development.

For these reasons, the present study chose Level 2 word-families rather than Level 6 word-families as the unit of counting for the EWL. However, unlike the traditional definition of lemmas which separate parts of speech, the present study defined Level 2 families as a word-form (headword) plus its inflections without distinguishing between parts of speech. This expanded version of lemmas have been called flemmas (family lemmas), but in this study we will refer to them as Level 2 families.

**Research questions**

The aim of the present study is to develop a wordlist for L2 beginners by including the best items in terms of lexical coverage from the GSL, BNC2000, BNC/COCA2000, and New-GSL. It sought to develop the Essential Word List (EWL) through answering the following seven questions.

1. What is the mean coverage provided by each set of 100 Level 2 headwords from a master list made up of Level 2 headwords from the GSL, BNC2000, BNC/COCA2000, and New-GSL in 18 corpora?
2. What is the mean coverage provided by each set of 100 Level 2 headwords plus members from the master list in the 18 corpora?
3. How many headwords should be included in an EWL?
4. Do the EWL headwords provide higher mean coverage in 18 corpora than the best headwords from each of the source lists from which the EWL was developed (GSL, BNC2000, BNC/COCA2000, and New-GSL)?
5. Do the EWL families provide higher mean coverage in 18 corpora than the best families from each of the source lists?
6. What is the overlap between the EWL headwords and the best headwords from the master list that were found in nine spoken corpora?
7. What is the overlap between the EWL headwords and the best headwords from the master list that were found in nine written corpora?

**Materials**

*The master list*

A master list was created of Level 2 word family (flemma) headwords from four source lists: West's (1953) GSL, Nation's (2006) BNC2000, Nation's (2012) BNC/COCA2000, and Brezina and Gablasova's

(2013) New-GSL. Because word families at Level 2 of Bauer and Nation (1993) were chosen as the unit of counting in the present study while word-families at Level 6 were the unit of counting in the original versions of the GSL, BNC2000, and BNC/COCA2000, Level 6 word-families from these lists were converted into Level 2 families. This was done by regrouping the GSL, BNC 2000, and BNC/COCA 2000 word-family members by following Leech, Rayson, and Wilson's (2001) principles for creating lemmatized wordlists. For example, the word-family *study* has six members: *study, studied, studies, studying*, *studious,* and *studiously.* When converted into flemmas, these members were grouped into three families: *study* (*study, studied, studies*, *studying*), *studious* (*studious),* and *studiously* (*studiously*). Once the conversion had been completed, the Level 2 word family versions of the GSL, BNC2000 and BNC/COCA2000 had 6,601, 6,465, and 6,412 headwords, respectively.

Because there was overlap between the items in the four lists, repeated headwords were excluded, resulting in 8,722 headwords remaining in the master list. A further 66 headwords were excluded. These items were letters (e.g., *B, X)* (20), affixes (e.g., *anti*, *non*) (7), cities (2), people's names (2), and the names of places and languages (35). Although learning letters of the alphabet is important, letters were excluded because it was assumed that L2 beginners would know them before learning English words. Learning affixes, especially the high-frequency affixes, has value for L2 beginners because they may have insufficient English morphological knowledge. However, it may be more reasonable to introduce a list of affixes when learners have reached a certain level rather than introducing affixes together with words right at the beginning (Nation, 2013). Proper nouns such as cities' names and people's names were not included because they are usually transparent and may have less value to learners than content words. The names of places and languages were excluded for two reasons. First, to be consistent with the decision to exclude the names of people and cities, these proper nouns should not be included in the master list. Second, the 35 names of places and languages may be biased towards the corpora from which the source lists were developed. For example, *Scot*, a BNC headword, appeared 496 times in the BNC but not in several corpora of other English varieties. This suggests that *Scot* was included in the BNC2000 not because it is a high-frequency word, but because it occurred very frequently in the BNC, the corpus from which the BNC2000 was developed.

Table 15.1: Nine spoken corpora used in the present study

| Name | Tokens | Variety of English |
|---|---|---|
| Spoken corpora | | |
| British National Corpus (spoken component) | 10,484,320 | British |
| International Corpus of English (spoken component) | 5,641,642 | Indian, Philippino, Singapore, Canadian, Hong Kong, Irish, Jamaican & New Zealand |
| Open American National Corpus (spoken component) | 3,243,449 | American |
| Webb and Rodgers (2009a) movie corpus | 2,841,573 | British & American |
| Wellington Corpus of Spoken New Zealand-English | 1,112,905 | New Zealand |
| Hong Kong Corpus of Spoken English | 977,923 | Hong Kong |
| Webb and Rodgers (2009b) TV program corpora | 943,110 | British & American |
| London-Lund corpus | 512,801 | British |
| Santa Barbara Corpus of Spoken American-English | 320,496 | American |

Table 15.2: Nine written corpora used in the present study

| Name | Tokens | Variety of English |
|---|---|---|
| British National Corpus (written component) | 87,602,389 | British |
| Open American National Corpus(written component) | 12,839,527 | American |
| International Corpus of English(spoken component) | 3,467,451 | Indian, Philippino, Singapore, Canadian, Hong Kong, Irish, Jamaican, New Zealand & American |
| Freiburg-Brown corpus of American-English | 1,024,320 | American |
| Freiburg–LOB Corpus of British-English | 1,021,357 | British |
| Wellington Corpus of Written New Zealand- English | 1,019,642 | New Zealand |
| Lancaster-Oslo/Bergen corpus | 1,018,455 | British |
| Brown corpus | 1,017,502 | American |
| Kolhapur Corpus of Indian-English | 1,011,760 | Indian |

### *The corpora*

Nine spoken and nine written corpora were used in the present study to examine the coverage provided by the headwords from the master list (Tables 15.1 and 15.2). These 18 corpora were in the form of untagged text files.  They varied in terms of size, type of discourse, and variety of English. The number of tokens ranged from 320,496 to 10,484,320 in the spoken corpora, and from 1,011,760 to 87,602,389 in the written corpora. The corpora represented 10 varieties of English: American-English, British-English, Canadian-English, Hong Kong-English, Indian-English, Irish-English, Jamaican-English, New Zealand-English, Philippino-English, and Singapore-English. Thus, it was expected that the 18 corpora would provide a thorough picture of the vocabulary that is essential for L2 beginners.

### Procedure

This study had three phases: (1) ranking the Level 2 headwords in the master list according to the mean coverage they provided in the 18 corpora, (2) determining the number of headwords to include in the EWL, and (3) assessing the EWL. Phase 1 related to determining the relative value of items in the four source lists. Phase 2 determined the cut-off point of the EWL. Phase 3 focused on evaluating the EWL.

### *Ranking the headwords in the master list*

Four steps were followed to determine the ranking of the headwords in the master list. First, the frequency of each headword was examined in each corpus. This was done by running each corpus through Nation, Heatley, and Coxhead's (2002) RANGE program with the master list in turn serving as the baseword list. RANGE is a program which analyses the lexical coverage provided by a wordlist in a text. It can be downloaded from Paul Nation's website (http://www.victoria.ac.nz/lals/about/staff/paul-nation). The second step was to calculate the coverage provided by each headword in each corpus. In this step, the frequency of each headword was divided by the number of running words in the corpus and multiplied by 100. For example, the coverage of *programme* in the Wellington Corpus of Spoken New Zealand-English (WSC) was 0.015% (165÷1,112,905 x 100 = 0.015%). The third step was to calculate the mean coverage of each Level 2 family in all 18 corpora. This was done by adding the coverage provided by the headwords in each of the 18 corpora and then dividing by the number of corpora (18). Mean coverage of the headwords in each corpus was more useful than the combined frequencies because combined frequency would bias the results towards findings in the largest corpora. By using the mean coverage of the headwords across 18 different corpora, range of lexical coverage was a key

criterion to rank the items in the present study. The fourth step was to rank headwords from the master list according to their mean coverage. That is, Level 2 word family headwords with the largest mean coverage were at the top of the master list while Level 2 word family headwords with the smallest mean coverage were at the bottom.

### Determining the number of EWL headwords

To determine how many headwords should be included in the EWL, two steps were followed. In the first step, the mean coverage provided by each set of 100 headwords from the master list and by each set of 100 Level 2 word families were determined. The present study examined the mean coverage provided by master list items at every 100-lemma headword level up to the 2,000-headword level. The mean coverage provided by each set of 100 headwords was calculated by adding the mean coverage of each headword in the set together. For example, the mean coverage provided by the set of the 1$^{st}$ 100-headwords was the sum of the mean coverage of each item in the top 100 headwords of the master list. To determine the mean coverage provided by sets of 100 Level 2 families, the coverage provided by each set of 100 families in each corpus was determined by running each corpus though RANGE with each set serving as the baseword list. Then, the mean coverage provided by each set of 100 Level 2 word families was calculated by adding the coverage in each of the 18 corpora together and dividing by 18.

In the second step, the cut-off point of the EWL was decided based on three criteria: practicability, change in the lexical coverage curve, and amount of lexical coverage. Practicability considered the size of the EWL in relation to the feasible amount of vocabulary that can be acquired by L2 leaners within a language program. The purpose of the present study is to develop a more practical wordlist for L2 beginners; therefore, practicability was the primary criterion to determine the length of the EWL. It would influence the decision related to the other two criteria: change in the lexical coverage curve, and amount of lexical coverage. Change in the lexical coverage curve involved examining the change in the lexical coverage provided by each set of 100 Level 2 headwords, and by these headwords plus members. Coverage by headwords is the actual coverage that learners may gain if they know the headwords. Coverage by headwords plus members reflects the potential coverage that learners may achieve if they can recognize members of these headwords. Although headwords were chosen as the primary unit of counting in the present study because they were usually the most frequent member in a lemma, there are still chances that members are more frequent than headwords. Therefore, it is also useful to use coverage provided by Level 2 word families as one criterion. Using both units of counting to decide the cut-off point provides an indication of how knowledge of two related but different units of counting might affect comprehension. Amount of lexical coverage examined the number of words necessary to reach different lexical coverage figures. Earlier studies have decided the length of a list based on the amount of vocabulary necessary to reach 95% coverage of text. However, lower coverage figures may still provide some indication of learners' progress in overall language development and assist teachers and course designers in organizing their English language programs to support learners' comprehension, as well as their lexical development. The number of words needed to reach different coverage figures was determined by the cumulative coverage provided by each set of 100 Level 2 headwords, and by the cumulative coverage provided by these headwords plus members.

### Assessing the EWL

Four criteria were used to evaluate the EWL. The first criterion involved a comparison between the mean coverage provided by the EWL headwords in the 18 corpora and the best headwords in terms of lexical coverage from the four source lists. The second criterion compared the mean coverage provided by the EWL word families with the mean coverage provided by the best word families from each source

list. The mean coverage provided by the best items in terms of lexical coverage in each source list was determined by following the same steps used to find the mean coverage provided by the EWL items. The third criterion was the overlap between the EWL headwords and the best headwords in terms of lexical coverage from the master list that were found in nine spoken corpora. The fourth criterion was the overlap between the EWL headwords and the best headwords in terms of lexical coverage from the master list that were found in nine written corpora. To determine the best headwords in nine spoken corpora, and in nine written corpora, the same steps used to select the EWL headwords were followed. Using both coverage provided by headwords and coverage provided by Level 2 word families as criteria to compare the EWL with the four source list provided a better picture about the actual coverage and potential coverage that L2 beginners may gain by knowing these wordlists. Looking at the mean coverage of the EWL and the best items in the source lists in the 18 corpora demonstrated the relative value of the lists in general, while the overlap between the EWL headwords and the best headwords in spoken and written corpora assessed the value of the EWL in different kinds of discourse. Together, these four criteria should provide a thorough assessment of the EWL.

Table 15.3: Additional coverage provided by the master list headwords and members at each 100 lemma headword level in 18 corpora

| Headword level | Examples | Coverage provided by each set of 100 words (%) | |
| --- | --- | --- | --- |
| | | Headwords | Headwords & members |
| $1^{st}$100 | the, okay | 45.68 | 55.46 |
| $2^{nd}$ 100 | sure, maybe | 5.62 | 6.71 |
| $3^{rd}$ 100 | sorry, hey | 2.94 | 3.71 |
| $4^{th}$ 100 | please, run | 2.08 | 2.70 |
| $5^{th}$ 100 | alright, hi | 1.61 | 2.06 |
| $6^{th}$ 100 | thanks, ok | 1.29 | 1.78 |
| $7^{th}$ 100 | hello, bye | 1.04 | 1.36 |
| $8^{th}$ 100 | drink, fast | 0.89 | 1.31 |
| $9^{th}$ 100 | tea, heavy | 0.77 | 1.11 |
| $10^{th}$ 100 | garden, huge | 0.69 | 0.99 |
| $11^{th}$ 100 | busy, weather | 0.62 | 0.90 |
| $12^{th}$ 100 | fresh, draw | 0.56 | 0.76 |
| $13^{th}$ 100 | active, holiday | 0.51 | 0.72 |
| $14^{th}$ 100 | fire, ride | 0.46 | 0.63 |
| $15^{th}$ 100 | shoot, lake | 0.41 | 0.61 |
| $16^{th}$ 100 | tiny, neck | 0.37 | 0.54 |
| $17^{th}$100 | vast, snow | 0.34 | 0.49 |
| $18^{th}$ 100 | attractive, channel | 0.32 | 0.45 |
| $19^{th}$100 | journey, calm | 0.29 | 0.43 |
| $20^{th}$ 100 | consumer, loud | 0.27 | 0.43 |

**Ranking the headwords in the master list**

The coverage provided by the sets of 100 headwords from the master list up to 2,000 headwords, as well as examples of items from each set are shown in Table 15.3. The mean coverage figures for the items in the different sets reflect their varying relative values. Those at higher levels are of greater value to language learners than those at lower levels. The 1ˢᵗ 100-headword level included items such as *the* and *okay*. The 10ᵗʰ 100-headword level included items such as *garden* and *huge*. The 20ᵗʰ 100-headword level included items such as *consumer* and *loud*.

In answer to Research Questions 1 and 2, the 1ˢᵗ 100-headwords provided mean coverage of 45.68% and the 1ˢᵗ 100 Level 2 word families (flemmas) provided 55.46% coverage. After the 1ˢᵗ 100-headwords, the mean coverage fell quickly. The 2ⁿᵈ 100-headwords provided mean coverage of only 5.62%; plus members, they provided 6.71% coverage. The coverage provided by headwords from the 3ʳᵈ, 4ᵗʰ, 5ᵗʰ, 6ᵗʰ, and 7ᵗʰ 100-headword levels was 2.94%, 2.08%, 1.61%, 1.29%, and 1.04%, respectively. The coverage provided by these headwords with their members was 3.71% (3ʳᵈ 100-headword level), 2.70% (4ᵗʰ 100-headword level), 2.06% (5ᵗʰ 100-headword level), 1.78% (6ᵗʰ 100-headword level), and 1.36% (7ᵗʰ 100-headword level). Beyond the 8ᵗʰ 100-headword level, the mean coverage provided by each 100-headword set was less than 1% while the mean coverage provided by families was less than 1% by the 10ᵗʰ 100-headword level.

**Determining the number of EWL headwords**

In answer to Research Question 3, the three criteria (practicability, change in the lexical coverage curve, and amount of lexical coverage) provide support for 800 items as the cut-off point for the Essential Word List. As the primary criterion, practicability will first be discussed alone, and then in relation to the other two criteria.

Practicability indicates that the EWL should have no more than 1,000 items. Earlier research on vocabulary growth has shown that L2 learners can acquire around 400 word families (Webb & Chang, 2012) or 500 lemmas (Milton, 2009) in a year. With this modest vocabulary growth rate, learning a list of more than 1,000 items may be too ambitious a goal for L2 beginners within an institution. This is supported by earlier research showing that EFL students from a range of contexts often fail to master the 1ˢᵗ 1,000 items despite a lengthy period of English instruction (Webb & Chang, 2012; Henriksen & Danelund, in press; Nurweni & Read, 1999; Quinn, 1968). A wordlist of less than 1,000 items is a more feasible task that might be learned within a single institution over two years. It focuses learners' attention on the most important items, which provide a much larger amount of lexical coverage than the subsequent 1,000 items (Dang & Webb, under review; Engels, 1968).

Practicability was then considered together with the other two criteria to determine a cut-off point within the 1ˢᵗ 1,000-headword level. Figure 15.1 illustrates changes in the coverage curves up to the 10ᵗʰ 100-headword level. The lower line presents the coverage provided by sets of headwords while the upper line presents the coverage provided by sets of headword plus members. In both cases, there was a decline in the coverage provided by each set of 100-items as the headwords became less frequent. There was a huge drop in coverage between the 1ˢᵗ and 2ⁿᵈ 100-headword level. From the 2ⁿᵈ to the 8ᵗʰ 100-headword levels, the amount of additional coverage, though not as high as that at the 1ˢᵗ 100-headword level, was still relatively large. However, beyond the 8ᵗʰ 100-headword level, the curve flattens out and the amount of additional coverage was less than 1%. The small change in the coverage between sets of 100 headwords beyond the 800 cut-off point suggests that the sequencing of items becomes less reliable because of the small difference in the mean coverage provided by headwords in

adjoining levels. That is, items which are in the 9[th]100 could just as well be in the 10[th]100. The lexical coverage curve criterion suggested two possible cut-off points for the EWL: 100 words and 800 words. If only the lexical coverage curve were used as the criterion to determine the cut-off point, 100 words would have been a more reasonable option because there is an extremely large decrease in coverage between the 1[st] and 2[nd] 100-headword sets. However, when the lexical coverage curve was considered together with practicability, 800 words is a better option. 78% of the words in the 800-item list were lexical words, and 22% were function words. In contrast, the percentage of lexical words and function words in the 100-item list was 28% and 72%, respectively. Lexical words are "words that convey content meaning" while function words are "words that express grammatical relationship" (Biber, Conrad, & Leech, 2002: 457-458). As lexical words enable L2 beginners to express their ideas, a list with an insufficient number of lexical words may not be very useful for L2 beginners. Therefore, an 800-item list seems more appropriate than a 100-item list when the coverage curves and practicability were considered together.



Figure 15.1. Coverage by each set of 100 headwords

Table 15.4: Cumulative mean coverage provided by the master list headwords and members at each cut-off point in 18 corpora.

| | Cumulative coverage at each 100 Level 2 headword point (%) | |
|---|---|---|
| Number of headwords | Headwords | Headwords & members |
| 100 | 45.68 | 55.46 |
| 200 | 51.3 | 62.17 |
| 300 | 54.24 | 65.88 |
| 400 | 56.32 | 68.58 |
| 500 | 57.93 | 70.64 |
| 600 | 59.22 | 72.42 |
| 700 | 60.26 | 73.78 |
| 800 | 61.15 | 75.09 |
| 900 | 61.92 | 76.2 |
| 1,000 | 62.61 | 77.19 |

The 800 item cut-off point was also supported by the third criterion (amount of lexical coverage). The top 800 headwords provided mean coverage of 61.15%, and potential coverage of 75.09% if all members of the lemmas were known (Table 15.4). The purpose of the EWL is to provide L2 beginners with the foundation for further vocabulary learning. Learning a relatively small number of words but reaching the 60% and 75% levels of coverage might be considered meaningful and practical to all stakeholders: teachers, program coordinators, and students. In this case, learning the 800 headwords would allow students to recognize over 60% of English words and as much as 75% of the English language if all members of the Level 2 families are known. The pedagogical significance of gaining knowledge of such a large proportion of English through studying a relatively short wordlist should be motivating to all stakeholders. Taken together, the three criteria suggested that 800 items should be the number of items in the EWL. The EWL is included in Appendix 3.

**Assessing the EWL**

In answer to Research Questions 4 and 5, the EWL headwords and families provided higher mean coverage in the 18 corpora than the best 800 headwords and similarly-sized families from each source list. The coverage provided by the EWL headwords in the 18 corpora was 61.16%. This is higher than the coverage provided by the top 800 GSL headwords (57.86%), top 800 BNC2000 headwords (57.66%), top 800 BNC/COCA2000 headwords (58.39%), and top 800 New-GSL headwords (60.83%). Similarly, the EWL Level 2 families provided higher mean coverage in 18 corpora (75.09%) than the best 800 Level 2 families from each source list (72.24%, 71.63%, 72.72%, 74.92%). This is not surprising because the EWL headwords were the best items from the four source lists. The fact that the EWL families provided the highest coverage also strongly supports the choice of headwords as the primary unit of counting in the present study.

It might be assumed that the top ranked 800 items in the best source list in the comparisons (New-GSL) are a reasonable substitute for the EWL. However, our analysis indicated that while there were many strong items in the New-GSL (and the other three lists), the rank order of the items is quite different when based on their coverage in the 18 corpora. There were 186 different items in the EWL and the top 800 items in the New-GSL indicating that the lists are quite different, and that the EWL is not simply a replica of the New-GSL.

In answer to Research Questions 6 and 7, 86.5% of the EWL headwords (692 items) appeared in the best 800 spoken headwords, and 698 EWL headwords (87.25%) were included in the best 800 written headwords. Importantly, 590 out of 800 EWL headwords (73.75%) appeared in both the best 800 spoken headwords and the best 800 written headwords. Among the 210 remaining headwords, 102 headwords (12.75%) appeared in the top 800 spoken headwords alone and 108 headwords (13.5%) appeared in the top 800 written headwords alone. The fact that most EWL headwords appeared in both the top 800 spoken and the top 800 written headwords, and there was a good balance in the number of remaining headwords that were unique to spoken and written discourse indicates that the EWL included basic words that are necessary for both written and spoken texts. This suggests that it would likely meet the needs of L2 beginners.

**Discussion**

As well as providing higher mean coverage in 18 corpora, the EWL has seven other strengths that make it superior to the source lists. First, unlike the four source lists, the number of items in the EWL was determined by examining the issue from different perspectives with the characteristics of target list-users (L2 beginners) in mind. Therefore, it may better reflect L2 beginners' needs.

Second, the EWL items may have greater validity than those from the four source lists. Unlike the four source lists, in the selection of EWL words, the frequencies of 110 words that can be either proper nouns or common words (e.g., *frank, mark*) were adjusted to reflect the real value of learning these items. That is, the frequency of headwords that occurred as proper nouns was subtracted from the total frequency of the headwords in the corpus. For example, in the WSC, *mark* appeared 176 times in total, but it was used as a proper noun 77 times. Therefore, the final frequency of *mark* in the WSC was 99. Without this adjustment, *mark* would be among the top 800 headwords of the master list. Also, the frequency of 92 headwords which had American variants was adjusted by adding the frequency of American variants to the total frequency. For instance, the final frequency of *programme* in the WSC (165) was the sum of the frequency of the British variant (*programme*) (148) and its American variant, (*program*) (17). Counting frequencies of both British and American variants in the final frequency of the headwords ensures that the EWL will better represent the essential vocabulary that learners often encounter in different language contexts. Moreover, while the other lists included letters (BNC2000), proper nouns (BNC2000), and affixes (BNC2000, BNC/COCA2000), the EWL excluded these items. Without this treatment, 12 names of places (e.g., *Indian, London*), 18 letters (e.g., *B, Y*), one affix (*non*), and seven items that can be either proper nouns or common words (e.g., *mark, lord*) would be included in the EWL. This would have meant excluding 39 items from the EWL including: *colour, dollar, fight, park,* and *television*. Compared with names of places, letters and affixes, these words should provide greater value to L2 beginners.

Third, the EWL included items which are very common in general conversation but are absent from some source lists. For example, *okay* and *alright* do not appear in the GSL and New-GSL; *hey, hi, hello,* and *bye* are absent from the New-GSL. A wordlist which contains common words in general spoken conversation might be more valuable for L2 beginners because, "for most people, the spoken language

is the main source of exposure to language, and is thus the main engine for language change and dynamism"(McCarthy & Carter, 1997:38). With widespread use of Communicative Language Teaching and Task-based Language Teaching approaches that pay more attention to spoken language, a list with a considerable number of words common in spoken discourse may be very attractive to teachers and learners.

Fourth, unlike the GSL, BNC2000, and BNC/COCA2000, which use Level 6 word-families as the unit of counting, the EWL uses Level 2 families (flemmas). This is a more reasonable decision because the EWL does not require sophisticated morphological knowledge or include low-frequency lemmas. Therefore, it is more appropriate for L2 beginners who are unlikely to be able to recognize many family members.

Fifth, the EWL items were derived from their lexical coverage in 18 corpora representing different discourse types, and 10 different varieties of English. In contrast, creation of the items in the earlier lists was based on a maximum of four corpora. Moreover, frequencies of both American and British variants were counted in the development of the EWL. Hence, the EWL should better represent the essential vocabulary encountered by learners in diverse situations.

Sixth, while none of the four source lists distinguish between function words (e.g., *the, of, in, at*) and lexical words (e.g., *know, big, people*), the EWL was divided into a list of 624 lexical words and a list of 176 function words. Although there are a number ways of classifying function words and lexical words, to be consistent, the present study follows Biber et al.'s (2002) classification. Words which can be either function words or lexical words (e.g., *have, past*) will be considered function words. However, to allow flexibility in the implementation of the EWL, teachers and learners can reclassify some EWL items into function word or lexical word lists. Classifying the EWL items into function words and lexical words has pedagogical value because of their different characteristics. In a text, lexical words are more salient than function words; therefore, 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. However, it is more reasonable to incorporate teaching function words with other components of language lessons due to their lack of salience in the text. No other word list has made the distinction between lexical and function words. This also makes the EWL more pedagogically appropriate.

Seventh, the EWL list of lexical words has sub-lists with manageable sizes. While the other lists either do not have sub-lists (New-GSL) or have 1,000-item sub-lists that might be too large to be incorporated effectively into language learning programs (GSL, BNC2000, BNC/COCA2000), the EWL list of lexical words is divided into 13 sub-lists according to decreasing mean coverage. The first 12 sub-lists have 50 headwords each while Sub-list 13 has 24 headwords. The mean coverage provided by each sub-list ranges from 6.26% (Sub-list 1) to 0.20% (Sub-list 13). Breaking the EWL list of lexical words into 50-headword sub-lists has two benefits. First, the size of the sub-lists is small enough to fit into individual courses within an English language program. Second, teaching the EWL lexical words following the rank order of sub-lists will increase learning effectiveness because it ensures that the most useful items are learned first. It also allows programs to prepare a curriculum that covers all sub-lists, and avoids teaching the same items between courses.

With these strengths, the EWL is a more suitable list for L2 beginners than the four source lists. Considering the influence of the GSL in vocabulary learning and practice, it is hoped that in the long run the EWL will receive the same attention from textbook authors, course designers, teachers, learners and researchers. However, promoting the use of the EWL does not mean that the present study does not recognize the value of the four source lists. The GSL, BNC2000, BNC/COCA2000, and New-GSL still have

value, but perhaps they are more useful for intermediate-level learners and researchers rather than L2 beginners.