# The Newspaper Word List: A Specialised Vocabulary for Reading Newspapers

## (Teresa) Mihwa Chung
*Korea University of Sejong, South Korea*

The primary purpose of this study is to identify words that are of special importance for reading newspapers. In the Newspaper Corpus of 579,849 running words, 588 word families are identified as Newspaper Words. These words account for 6.8% of the running words in the corpus. When combined, proper names and 2,521 families of the *General Service List of English Words* (GSL) and the NWL make up 92.5% of the running words in the corpus. This is lower than the 98% ideal coverage required for understanding a text successfully, but very high given the small vocabulary size. Thus, the NWL will give the best return for vocabulary learning to learners of English as a foreign language who wish to read newspapers as soon as possible.

　本研究の目的は新聞を読むのに必要な語彙を特定することである。Newspaperコーパスに記載された579,849語の中から6.8％に上る588語をNewspaper　Wordsとして選び出した。固有名詞、General Service List of English Wordsの2,521語、NWLの語彙を総計するとコーパスの92.5％になる。テキストを理解するのに必要とされている98％よりはやや低い数値であるが、NWL全体の総語数を考慮すれば非常に高い値であるということができる。したがってNWLは新聞英語を早く読めるようになりたいと望む英語学習者にとっては最も効率のよいものであるということができる。

R eading is one of the most common and important ways of learning another language and one major type of reading material is newspapers. Newspapers are often used in reading classes in order to develop reading skills and expand vocabulary knowledge (Hwang & Nation, 1989; Klinmanee & Sopprasong, 1997). There are several reasons for this. Firstly, newspapers are easily and cheaply available in hard copy or online. Secondly, newspapers are authentic materials that are commonly read by the native speakers of the language. Thirdly, they provide a wide choice of interesting topics from which teachers or learners can choose reading texts. Finally, reading newspapers is considered to be not only a good way of reviewing old vocabulary learned, but also of learning new vocabulary from context (Hwang & Nation, 1989).

Despite these potential advantages, many learners find it difficult to read unsimplified newspaper texts. There are a number of factors that contribute to difficulty in reading, but vocabulary knowledge has consistently been found to be the most influential factor affecting comprehension (Hirsh & Nation, 1992; Nation & Coady, 1988).

Hu and Nation (2000) suggested that knowledge of at least 98% of the total words (tokens or running words) in a text is the minimum required for adequate reading comprehension. A recent study by Nation (2006) examined the receptive vocabulary size needed for reading newspapers using the reportage category of the parallel LOB [the Lancaster Oslo/Bergen Corpus of British English] (Johansson, 1978), FLOB (Hundt, Sand, & Siemund, 1999), Brown (Francis & Kučera, 1979), and Frown [the Freiburg-Brown Corpus of American English] (Hundt, Sand, & Skandera, 1999) corpora. In that study, Nation used the British National Corpus (BNC) in order to develop word frequency lists and applied them to other corpora, including newspapers. He estimated that knowledge of the 8,000 most frequent word families and proper names is needed to reach 98% lexical coverage. This number represents a fairly large vocabulary, particularly for adult learners of English who want to read newspapers as a means of learning English and for knowing what is occurring inside and outside of the country. An important methodological approach used by Nation (2006) involved using a reportage news category without dividing it into smaller news sections (e.g., international news and business news); as a result, a range criterion was not included, which is important when selecting a wide range of words occurring with high frequency in newspapers. In the present study, four news sections are examined in order to obtain more detailed results from coherent sections.

It is important to remember that particular words common in certain kinds of writing occur frequently in those texts, and therefore provide good coverage for those text types. A good example of a specialised vocabulary is the Academic Word List (Coxhead, 2000), which consists of only 570 word families, and provides coverage of at least 9% of the running words in a wide range of academic texts. As another example, Ward (1999) notes that a vocabulary of only 2000 word families provides 95% coverage of the tokens in many engineering texts, which is sufficient for 1st-year students to read required textbooks. Learning such a specialized vocabulary provides learners with a shortcut to coping with the vocabulary of such texts.

For this reason, if researchers develop specialized vocabulary lists, only a small vocabulary is needed to make certain types of texts easily accessible provided the specialized list of vocabulary is acquired. In addition, when teachers narrow the focus on teaching vocabulary such as for reading newspapers and engineering texts, the vocabulary burden required of learners is lowered and as a result, learners benefit from such instruction. Therefore, well-chosen specialized vocabulary lists can give learners the best results with the least effort.

To date, it is not known how large this specialized newspaper vocabulary might be and what kinds of words would form it. Thus, the primary purpose of this study is to identify the specialized vocabulary of newspapers. Three research questions will be investigated:

1. How many word families make up a specialised vocabulary of newspaper texts (hereinafter, the Newspaper Word List)?
2. What percentage of the tokens in newspaper texts does the Newspaper Word List cover?
3. How often do the words in the Newspaper Word List occur in different newspapers and different news divisions?

## Materials and Methods

### Computer Programs

The analyses in this study were performed with the vocabulary analysis program *Range32H* (Heatley & Nation, 2006) in order to count and create a list of proper names and Newspaper Words. The program uses two Baseword lists: Baseword one is the first 1,000 words and Baseword two is the second 1,000 words of the GSL.

A weakness of *Range* is that the program cannot distinguish *ward* (a family name) from *ward* (a section of a hospital). This problem was addressed by looking at the context in which the word was embedded. Similarly, *aid* as a verb is not distinguished from *aid* as a noun; however, this kind of polysemic use was not considered a problem because the learning burden of the noun form is very small if the verb use is known. Many English verbs may commonly be also used as nouns.

### Determining a Unit of Counting Word Families

In this study, the word family is used as the unit when counting words. The level of word family used here is composed of a base form together with its inflected forms and derived forms as described in Level 6 of Bauer and Nation's scale (1993). A word family represents a group of words whose forms and meanings are closely related to each other and which can be understood with little or no extra learning when one or more of the members is already well known to a learner. Thus, word types from the same word family are counted as the same word. Both American and British spellings are counted in the same family. For example, *analyse* and *analyze* are counted in the family *analyse*. The main justification for the use of the word family is that it best represents the kind of knowledge needed when meeting words in reading, and the goal of this study is to examine the vocabulary needed for reading newspapers. Table 1 illustrates how large word families can be. Each word in italics is the most frequently occurring form in that family in the Newspaper Corpus.

### Compiling the Newspaper Corpus

The news texts used in this study were obtained from the Internet Public Library drawing on texts published from 23 February to 23 May 2006. All texts were obtained in electronic form. The dates of the reports and the names of the reporters and the newspapers were removed.

In making the Newspaper Corpus, four principles were followed. The first principle was that newspapers for the Newspaper Corpus of English (hereinafter, the Newspaper Corpus) had to represent the kinds of English newspapers that native speakers of English would typically read. Three newspapers were chosen: *The Dominion Post* from New Zealand, *The Independent* from the United Kingdom, and *The New York Times* from the United States of America. These are representatives of high quality English newspapers in these three English-speaking countries. Though the sensational tabloids

**Table 1. Examples of Three Word Families in the Newspaper Corpus**

| FINANCE | SECURE | INVEST |
|---|---|---|
| FINANCES | SECURES | INVESTS |
| FINANCED | SECURED | INVESTED |
| FINANCING | SECURING | INVESTING |
| *FINANCIAL* | SECURELY | *INVESTMENT* |
| FINANCIALLY | *SECURITY* | INVESTMENTS |
| FINANCIER | SECURITIES | INVESTOR |
| FINANCIERS | UNSECURED | INVESTORS |
| | INSECURE | REINVEST |
| | INSECURITY | REINVESTS |
| | INSECURITIES | REINVESTED |
| | | REINVESTING |
| | | REINVESTMENT |

are widely read, they were not selected because the Newspaper Word List is intended to help learners study English in an intensive course, often with the goal of going to university.

The second principle was that the corpus had to be large enough in order to allow the lower frequency candidates for a specialized vocabulary of newspapers to have a reasonable number of occurrences (Kennedy, 1998; Leech, 1987; Sinclair, 1991). A corpus of 579,849 running words proved to be large enough to obtain a minimum frequency of at least 20 occurrences of each candidate word.

The third principle was that the corpus had to contain approximately equal-sized, representative sections of each newspaper in order to measure the range of occurrence of words. Range is vital because lexical items that will be met when reading different sections of a newspaper and different newspapers should be selected.

The Newspaper Corpus consisted of 12 sections, namely the four main news divisions (Business, International, National, and Sports) from three newspapers (i.e., *The Dominion Post*, *The Independent* and *The New York Times*). Table 2 provides data concerning the size of the 12 news sections each counted by the *Range* program.

### Table 2. Tokens in each of the 12 News Sections

| News division | *The Dominion Post* | *The Independent* | *The New York Times* | Total |
|---|---|---|---|---|
| National | 48,270 | 47,816 | 48,527 | 144,613 |
| Business | 47,361 | 47,922 | 48,549 | 143,832 |
| Sports | 48,827 | 49,020 | 48,750 | 146,597 |
| International | 48,594 | 47,848 | 48,365 | 144,807 |
| Total | 193,052 | 192,606 | 194,191 | 579,849 |

As shown in Table 2, the National news texts in *The Dominion Post* contained 48,270 tokens and the combined National news texts from the three newspapers totaled 144,613 tokens. The four sections in the Dominion Post contained a total of 193,052 tokens. On average, each of the 12 news sections contained 48,300 tokens, each of the four news divisions 144,900 tokens, and each of the three newspapers a total of 193,000 tokens. Each of the 12 sections was of roughly equal size in order to obtain comparable statistical data from the various sections, and accordingly the frequency of the words was not biased by the size of each section (Leech, 1987; Sinclair, 1991).

The fourth principle was that texts in the corpus should be representative of news text types. Three conditions were considered. First, texts for the corpus should be selected from a news reportage category rather than from editorials, book and movie reviews, or advertisements because reporting news is a more typical function of a newspaper. Second, a large variety of news texts written by a large number of reporters should be included in the corpus. Third, the texts should be whole texts rather than a collection of partial texts, and relatively long texts need to be chosen in order to obtain specialized words with a higher frequency. Sampling whole texts gives topic-related words more opportunity to occur, though marked differences of individual writing style or topic might appear (Sinclair, 1991). Accordingly, the 868 texts comprising the Newspaper Corpus (see Table 3) were whole texts from reportage and 844 texts (97.2%) were between 200 and 2,000 words long; the shortest text was 131 words long and the longest was 5,054 words. A balance between short and long texts, and a balance in size between different news divisions were achieved where possible as shown in Table 3.

**Table 3. Number of Texts in Each News Division**

| News division | Number of texts |
|---|---|
| National | 221 |
| Business | 211 |
| Sports | 215 |
| International | 221 |
| Total | 868 |

Each news division contained 217 texts on average. Care was taken when compiling the corpus to ensure that texts were not repeated in different newspapers. There is a high risk of this occurring because wire services like Reuters and API provide news to newspapers all over the world. A very large amount of work was involved in collecting the corpus, as each of the 868 texts had to be downloaded one by one, checked to avoid duplication of texts, and edited for misspellings, spelling variations, foreign words, and hyphenated words (for details, see the section *Editing the News Texts* below).

### Editing the News Texts

The news texts were edited to make them computer readable and to avoid counting problems. After that, the texts were saved in the *plain text* format in order to make them suitable for analysis by the *Range* program.

*Hyphenated Words*: For hyphenated words with a deducible meaning from the meaning of each constituent (e.g., *large-scale*, *wide-bodied*, and *anti-war*), a space on each side of the hyphen was inserted using the Find and Replace function on the computer. This is because in terms of counting the occurrences of each word and measuring the vocabulary load of the text, it is better to count each constituent separately, as this avoids inflating the number of different word types. Where it is better to keep a hyphen in order to maintain the meaning of the whole, a hyphenated word was changed into one lexical item without a hyphen (e.g., *preemptive*, *email*, and *hiphop*).

*Foreign Letters*: When foreign words in *Word* format were saved in *plain text* format, the pre-existing *Word Document* format was lost, and this created problems in counting words. For example, *Löffler* was initially counted wrongly as two items, *L?* and *ffler*. **For this reason, foreign letters, for exam-**

ple *ö*, *é* and *á* as in *Löffler*, *René,* and *Chávez,* were replaced with the English letters *o*, *e* and *a*, respectively.

*Various Word Forms with the Same Meaning*: In the case of varying word forms with the same meaning such as *per cent* and *percent*, *per cent* was replaced with *percent*. Otherwise, *per cent* and *percent* would be counted as three items, *per*, *cent* and *percent*.

*Names with an Apostrophe*: Words written with an apostrophe, such as *Shi'ite* and *Fa'atau,* were rewritten as *Shiite* and *Faatau* in order to avoid each being counted as two items.

### Setting up Criteria for Identifying Specialized Words for Reading Newspapers

Three criteria were set up in order to ensure that the words identified were specialized vocabulary for reading newspapers.

*Special Purpose Vocabulary*: The first criterion was that newspaper words must be special purpose vocabulary. This meant that they could not be part of the high-frequency 2,000 words of English as defined by West's (1953) *General Service List of English Words* (GSL). In addition, no proper names were included on the list. One reason for choosing the GSL was to make the data comparable with the Academic Word List which also assumes knowledge of the GSL.

*Wide Range*: The second criterion of range had the highest priority because words should occur in a wide range of different news texts. In this study, range was measured by (a) determining the number of news divisions in which each candidate word occurred and (b) by counting the number of news sections across the three newspapers and the four news divisions in which the word was found (e.g., *The Dominion Post's* National news section and *The Independent's National*'s news section). Thus, Newspaper Words must occur in all four news divisions of the corpus, and 6 or more of the 12 smaller news sections. Because the primary aim of the study is to create a list of the most useful Newspaper Words rather than create a complete list of Newspaper Words, a range of 6 or more out of 12 was considered sufficient for identifying Newspaper Words.

*High frequency*: The third criterion was the total frequency with which the candidate words occurred in the Newspaper Corpus. Frequency is important but not foremost because creating a word list based on frequency alone allows a bias towards longer texts and topic-related words. In this study, Newspaper Words must occur 20 times or more in at least 6 out of the 12 sections in the corpus. The frequency cutoff point of 20 occurrences was chosen because in terms of practicality, 20 examples provide enough examples for a useful concordance analysis of an item (Leech, 1987; Sinclair, 1991).

### Making a List of Proper Names

In order to prevent frequently occurring proper names from being selected as Newspaper Words, a list of proper names was created by examining the words with a frequency of 20 or more occurring outside the GSL 2,000 words. The list of proper names included personal names (e.g., *Mary* and *David*), country names (e.g., *New Zealand* and *Britain*) and organization names (e.g., *Delta Air Lines* and *Duke University*). Abbreviations, such as NZQA, EU, and FIFA, were generally included in the list of proper names.

Certain items, such as *hawk* (as in Black Hawk helicopter and a kind of bird), *ward* (as in *Martin Ward* and a kind of room**),** *mount* (as in *Mount Tambora* and go up), and *range* (as in *Tararua Range* and widespread), were used as both a proper name and an ordinary item in the corpus. Items occurring more frequently as a proper name than an ordinary item in the Newspaper Corpus (e.g., *Hawk* and *Ward*) were placed in the list of proper names.

In order to make a list of proper names, all word types with a frequency of 20 or more that did not occur in the GSL were examined and a decision was made about which would be put into the list of proper names.

Note that after making a list of proper names, there are three Baseword lists to run with the *Range* program: Baseword one and two are the first 1,000 words and the second 1,000 words of the GSL; Baseword three is a list of proper names.

### Creating a List of Specialized Words for Reading Newspapers

The following steps were taken in order to identify all the word types outside the 2,000 words of the GSL and the list of proper names, to decide whether they met the criteria for identifying specialized words and thus to select potential candidates for a list of Newspaper Words.

First, word types occurring outside the three Baseword lists were identified by running the four news divisions of the Newspaper Corpus through the *Range* program. Second, 1,012 word types occurring in all four news divisions (Business, National, Sports, and International) and not in the GSL or proper name list were identified. Third, the 1,012 types were organized into 733 word families using the *Copy* function in the *Range* program drawing on Nation's fourteen 1,000-word lists from the British National Corpus. These lists have been carefully created and are a reliable source of word families. Fourth, by running the 12 news sections (see Table 2) of the corpus through the *Range* program, 523 word families with a range of 6 or more out of 12 and a frequency of 20 or more occurrences were identified.

Finally, word types occurring in only two or three news divisions were examined in order to determine whether counting word families rather than word types would allow more words to become candidates for the Newspaper Word List. This resulted in 65 word families (e.g., *adequate*, *bonus*, *consult*, and *score*) being added to the list, giving a total of 588 word families.

## Results

### *The Newspaper Word List and its Text Coverage*

From a corpus of 579,849 tokens, 588 word families (Appendix 1) were identified as specialized words for reading newspapers using the criteria of range and frequency. Table 4 shows how much of the Newspaper Corpus was covered by the GSL lists and the Newspaper Word List, and how many families in each list occurred in the corpus.

**Table 4. Text Coverage and Number of Families in Each List**

| Word list (Number of families in the list) | Coverage of the Newspaper Corpus | Number of families occurring in the Newspaper Corpus |
|---|---|---|
| Newspaper Word List (588 families) | 6.8% | 588 |
| Second 1,000 GSL (991 families) | 5.5% | 937 |
| First 1,000 GSL (998 families) | 74.2% | 996 |
| Total (2,577 families) | 86.5% | 2,521 |

Table 4 shows that the NWL covered 6.8% of the tokens in the corpus. This is higher than the 5.5% coverage of the second 1,000 GSL of the corpus. This contrast is even more striking when we consider that the total number of word families in the NWL (588 families) is much smaller than the 937 families occurring in the second 1,000 of the GSL.

The first 2,000 words of the GSL and the 588 newspaper word families in the corpus provide coverage of 86.5% of the running words in the corpus. This is a high degree of coverage with a relatively small number of words.

### The NWL Coverage of National, Business, Sports, and International News Divisions

Table 5 shows a comparison of the coverage of the four news divisions by the NWL, the GSL, and proper names.

**Table 5. Text Coverage of the Four News Divisions by Each List**

| Coverage | National news | Business news | Sports news | International news |
|---|---|---|---|---|
| NWL | 6.7% | 8.3% | 5.1% | 7.1% |
| Second 1,000 GSL | 5.9% | 5.5% | 5.5% | 5.3% |
| First 1,000 GSL | 74.8% | 74.7% | 74.9% | 72.5% |
| Proper names | 4.9% | 4.7% | 7.4% | 7.0% |
| Total coverage | 92.3% | 93.2% | 92.9% | 91.9% |

*Note*: Proper names are treated as known words because proper names are easily understood from the context or are already known to students (Hwang & Nation, 1989).

The NWL coverage of the Business news division is the highest (8.3%) and the coverage of the Sports news division is the lowest (5.1%). A factor contributing to the high coverage by the NWL of the Business news division is that some word families occurred extremely frequently in the Business news division, but were of much lower frequency in the Sports news as seen in Table 6.

The NWL coverage of the National and International news divisions is similar (6.7% and 7.1%, respectively). Within the most frequent top 10 words, three word families: *percent*, *issue*, and *secure* (*security* is the most frequent type) were included in both National and International news. The

other 7 families out of the top 10 included *labour*, *fund*, *drug*, *job*, *sex*, *investigation*, and *port* in the National news; and *military*, *protest*, *bomb*, *prime*, *terror*, *major*, and *region* in the International news.

**Table 6. A Comparison of the Number of Occurrences of 18 Word Families in the Business News and Sports News Sections**

| Word families | Number of Occurrences in Business News | Number of Occurrences in Sports News |
|---|---|---|
| PERCENT | 707 | 29 |
| INVEST | 456 | 9 |
| FINANCE | 195 | 18 |
| EXECUTIVE | 184 | 33 |
| FUND | 178 | 13 |
| ENERGY | 132 | 4 |
| CONSUME | 124 | 1 |
| REGULATE | 122 | 8 |
| BID | 112 | 32 |
| COMMISSION | 106 | 12 |
| ISSUE | 95 | 33 |
| PENSION | 93 | 1 |
| ANALYSE | 91 | 13 |
| CORPORATE | 84 | 3 |
| SHAREHOLDER | 81 | 1 |
| INCOME | 80 | 2 |
| EXPORT | 76 | 2 |
| REVENUE | 75 | 6 |

The second 1,000 words of the GSL had very similar coverage (about 5.5%) in all four news divisions, but the first 1,000 words of the International news division had slightly lower coverage (72.5%) than the others (around 74.8%).

The proper name coverage of the Sports news division was the highest (7.4%), while the Business news coverage was the lowest (4.7%). Because players' skills, team performances, and new players' names are mentioned frequently in the Sports News section, proper names occurred more frequently in this section than in any other news division. For this reason, if

the Sports news division is excluded from the Newspaper Corpus, the NWL coverage of the remaining combined texts rises to 7.4%. The smallest coverage by the NWL of the Sports news division was balanced by the biggest coverage of proper names. Thus, the total coverage of each of the four news divisions by the four combined lists was very similar, between 91.9% and 93.2%.

### The NWL Coverage of each Newspaper of the Three Countries

The coverage of the three newspaper corpora provided by the three word lists is shown in Table 7.

**Table 7. Coverage of Each Newspaper by the Three Word Lists**

|  | *The Dominion Post* | *The Independent* | *The New York Times* |
|---|---|---|---|
| NWL coverage | 7.3% | 6.5% | 6.6% |
| Second 1,000 GSL coverage | 5.7% | 5.6% | 5.2% |
| First 1,000 GSL coverage | 73.3% | 74.8% | 74.5% |
| Total | 86.3% | 86.9% | 86.3% |

As shown in Table 7, there is little difference in coverage provided by all three lists for the three newspapers. The NWL is clearly an international list and it could be expected to work well with other similar newspapers. Eight word families: *executive*, *final*, *invest*, *issue*, *major*, *percent*, *secure*, and *team*, were within the most frequent 20 words in the three newspapers.

### Range of the 588 Families of the NWL

Table 8 shows how many of the 588 families of the NWL occurred in 6 or more of the 12 sections of the Newspaper Corpus.

One hundred thirty-eight (24%) of the 588 word families occurred in all 12 news sections, and 567 families (96%) occurred in 7 or more sections. The wide range of the 588 families indicates that the list is likely to apply well to other similar quality newspapers.

### Table 8. Cumulative Number and Their Percent of 588 Word Families in Sections 6 to 12

| Number of news sections | Number of NWL families | Cumulative number | Cumulative per-cent of families |
|:---:|:---:|:---:|:---:|
| 12 | 138 | 138 | 24% |
| 11 | 127 | 265 | 45% |
| 10 | 110 | 375 | 64% |
| 9 | 94 | 469 | 80% |
| 8 | 60 | 529 | 90% |
| 7 | 38 | 567 | 96% |
| 6 | 21 | 588 | 100% |

### *Evaluation of the Newspaper Word List*

A frequency-based word list made from a particular corpus will provide reasonably high coverage of that corpus. In order to test whether the NWL provides good coverage of a different newspaper corpus, the newspaper sections of the Frown Corpus and the FLOB Corpus, both containing material written in the early 1990s, were chosen. These are relatively new compared with similarly structured but older Brown and LOB corpora compiled over 30 years earlier.

From the Frown and FLOB corpora, three categories (reportage, editorials, and reviews) were selected for making three corpora to test the NWL: (a) a reportage news corpus, (b) a reportage and editorials combined corpus, and (c) a reportage, editorials, and reviews combined corpus. The reportage news corpus contains 88 texts, amounting to 180,170 tokens; the reportage and editorials combined corpus, 142 texts (292,048 tokens); the reportage, editorials, and reviews combined corpus, 176 texts (362,584 tokens). Note that all texts of the Newspaper Corpus used in this study would be classified as reportage in the Frown and FLOB corpora. Table 9 shows a comparison of the three news corpora from the Frown and FLOB.

**Table 9. Number of NWL Families and Their Coverage in Various Corpora**

| Corpus (Tokens) | NWL Families | NWL Coverage |
|---|---|---|
| Reportage news (180,170 tokens) | 577 | 6.0% |
| Reportage and editorial combined (292,048 tokens) | 582 | 6.0% |
| Reportage, editorial, and reviews combined (362,584 tokens) | 582 | 5.7% |

As shown in Table 9, the coverage by the NWL of each of the three news corpora was similar at 6.0%, 6.0%, and 5.7%. This indicates that the NWL also works well with editorials and reviews sections. The 6.0% is slightly lower than the 6.8% coverage of the Newspaper Corpus compiled for this study. Five hundred seventy-seven out of the 588 newspaper families occurred in the reportage news texts, quite a lot given the smaller corpus size. The 11 NWL families which did not occur in the reportage corpus were *cellphone*, *detention*, *email*, *enrich* (*enrichment* is the most frequent type), *enroll*, *flu*, *immigrate* (*immigration* is the most frequent type), *internet*, *refine* (*refinery* is the most frequent word type), *virus*, and *website*. Such items are likely to be affected by the age of the corpus because there is more than a 15 year difference in age between the Newspaper corpus and the Frown and FLOB combined corpus. The items occurring frequently are also affected by new or emerging topics such as email and bird flu. Six families out of the 11 did not occur in the reportage and editorial combined news: *cellphone*, *email*, *enrich*, *flu*, *internet*, and *website*.

Newspaper texts from the Brown and LOB corpora were also examined to determine how much of the text the NWL covered. Because the Brown and LOB corpora were compiled in the 1960s (almost 50 years ago), the coverage of the NWL was around 5.1%, suggesting that the NWL is affected by current issues and needs to be updated periodically.

## Conclusions

In the Newspaper Corpus of 579,849 tokens, 588 word families were classified as Newspaper Words. The list of 588 families is a specialized vocabulary which provides a high coverage of newspaper texts. It accounted for 6.8% of the tokens in the corpus. One strength of the Newspaper Word List

is that the 588 families are a much smaller group than 937 families occurring in the second 1,000 GSL, but the coverage of the NWL is 1.3% better than the coverage by the second 1,000 GSL.

When combining the coverage of the NWL, GSL, and proper names, the coverage of the corpus comes to 93% coverage. Though this is lower than the 98% target coverage criterion specified by Hu and Nation (2000), this is very high and thus the NWL can provide second language learners who want to read English newspapers with a way to focus their vocabulary studies.

The NWL can add to the number of high frequency words that could be directly taught in class time and that deserve deliberate study by learners. It is important to remember that vocabulary learning should take place in a balance of activities, covering not only meaning-focused activities but also language-focused and fluency development activities. For maximum benefit, learners should read more related stories than unrelated stories. Following the same story through the several issues of the newspaper is an effective way of helping learners review the vocabulary learned previously (Hwang & Nation, 1989; Schmitt & Carter, 2000). **The NWL would be useful for teach**ers of English for specific purposes (ESP) who are interested in designing a vocabulary course for foreign language learners who wish to read newspapers as soon as possible.

In future studies, firstly, it may be desirable to collect data for more than 3 months and compile a bigger corpus covering a wider range of sections so that the NWL could be more widely applied in each newspaper. Secondly, the NWL may need to be updated every 5 to 7 years as it is partly influenced by current world events.

## Acknowledgements

*Mihwa Chung* teaches at Korea University of Sejong in Korea. Her doctoral thesis examined a range of ways of distinguishing technical terms from other words in English for Specific Purposes. Her current teaching and research interests include teaching and learning vocabulary (in particular, technical terms and specialized vocabularies), corpus analysis, reading courses and speed reading courses.

# References

Bauer, L., & Nation, P. (1993). Word families. *International Journal of Lexicography, 6*(4), 253-279.

Coxhead, A. (2000). A new academic word list. *TESOL Quarterly, 34*(2), 213-238.

Francis, W. N., & Kučera, H. (1979). *Manual of information to accompany a standard corpus of present day edited American English, for use with digital computers.* Providence, RI: Brown University.

Heatley, A., & Nation, P. (2006). Range32H (computer software). Wellington, New Zealand: Victoria University of Wellington.

Hirsh, D., & Nation, P. (1992). What vocabulary size is needed to read unsimplified texts for pleasure? *Reading in a Foreign Language, 8*(2), 689-696.

Hu, M., & Nation, P. (2000). Unknown vocabulary density and reading comprehension. *Reading in a Foreign Language, 13*(1), 403-430*.

Hundt, M., Sand, A., & Siemund, R. (1999). *Manual of information to accompany the Freiburg-LOB Corpus of British English (FLOB).* Freiburg: Englisches Seminar, Albert-Ludwigs-Universität Freiburg.

Hundt, M., Sand, A., & Skandera, P. (1999). *Manual of information to accompany the Freiburg-Brown Corpus of American English (Frown).* Freiburg: Englisches Seminar, Albert-Ludwigs-Universität Freiburg.

Hwang K., & Nation, P. (1989). Reducing the vocabulary load and encouraging vocabulary learning through reading newspapers. *Reading in a Foreign Language*, *6*(1), 323-335.

Johansson, S. (1978). *Manual of information to accompany the Lancaster-Oslo/Bergen Corpus of British English, for use with digital computers.* Oslo: University of Oslo.

Kennedy, G. (1998). *An introduction to corpus linguistics.* London: Longman.

Klinmanee, N., & Sopprasong, L. (1997). Bridging the vocabulary gap between secondary school and university: A Thai case study. *Guidelines*, *19*(1), 1-10.

Leech, G. (1987). General introduction. In R. Garside, G. Leech, & G. Sampson (Eds.), *The computational analysis of English: A corpus-based approach* (pp. 1-15)*. London: Longman.

Nation, P. (2006). How large a vocabulary is needed for reading and listening? *Canadian Modern Language Review*, *63*(1), 59-82.

Nation, P., & Coady, J. (1988). Vocabulary and reading. In R. Carter & M. McCarthy (Eds.), *Vocabulary and language teaching* (pp. 97-110). London: Longman.

Schmitt, N., & Carter, R. (2000). The lexical advantages of narrow reading for second language learners. *TESOL Journal*, *9*(1), 4-9.

Sinclair, J. (1991). *Corpus, concordance, collocation*. Oxford: Oxford University Press.

Ward**,** J. (1999). How large a vocabulary do EAP engineering students need? *Reading in a Foreign Language, 12*(2), 309-323.

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

## Appendix 1: Newspaper Words in 10 Sublists

The NWL is grouped into 10 sublists, and each sublist includes 60 words, except for Sublist 10 which includes 48. Groups of 60 break the learning task into units of a manageable size for a short-term course. The grouping of the sublists is based on range and frequency, and range is given precedence over frequency. Sublist 1 contains words which are of the widest range and the words in Sublist 10 are of the narrowest range among the 10 sublists.

The most frequently occurring member of each word family in the NWL is displayed in the list. Figures indicate the sublist of the NWL. For example, *abandoned* is the most frequent type of the members of the family *abandon*, and this family occurs in Sublist 8 of the NWL. Note that both American and British spellings are included in the word families (e.g., both *rumors* and *rumours* are included in the family *rumor*). Three prefixes (*pre-*, *ex-*, *pro-*) are included in the list because they are frequently used to make words, are predictable in meanings as in *pre-season, pre trial, pre-match, ex-adviser, ex-offenders, ex-employee, pro-democracy campaign, pro-Palestinian,* and *pro-life groups*, and are within the Level 6 affixes (Bauer & Nation, 1993) used for making a list of Newspaper words in this study.

### *Headwords of the Newspaper Word List in 10 Sublists*

| | | | | | |
|---|---|---|---|---|---|
| abandoned | 8 | adjusted | 3 | allegations | 5 |
| abuse | 5 | administration | 3 | alliance | 5 |
| academic | 3 | affected | 2 | allies | 4 |
| access | 3 | agenda | 3 | alongside | 5 |
| accompanied | 6 | aggressive | 2 | alternative | 4 |
| accurate | 6 | aid | 3 | amazing | 8 |
| achieve | 1 | airline | 8 | amid | 3 |
| acknowledged | 2 | airport | 4 | analysts | 1 |
| acquired | 2 | alarm | 5 | announced | 1 |
| adequate | 9 | alcohol | 8 | annual | 3 |

| | | | | | |
|---|---|---|---|---|---|
| anticipated | 5 | boost | 7 | cited | 4 |
| apartment | 7 | boss | 3 | clash | 7 |
| apparently | 2 | bounce | 7 | classic | 10 |
| appeal | 3 | brand | 5 | climate | 6 |
| approach | 1 | breach | 4 | clinical | 7 |
| appropriate | 3 | brewers | 8 | collapse | 4 |
| area | 1 | brief | 5 | colleagues | 7 |
| aspects | 7 | broker | 9 | column | 9 |
| assembly | 6 | budget | 3 | combat | 9 |
| asserts | 9 | bullet | 8 | comment | 1 |
| assessment | 2 | burden | 9 | commission | 3 |
| assets | 7 | bureau | 10 | committed | 1 |
| assistant | 2 | cabinet | 7 | communications | 6 |
| assume | 3 | campaign | 1 | community | 1 |
| assured | 2 | cancer | 10 | compensation | 6 |
| athletic | 10 | candidate | 3 | complex | 6 |
| attached | 3 | capable | 3 | comply | 5 |
| attitude | 7 | capacity | 7 | compound | 10 |
| attorney | 9 | captured | 5 | computer | 6 |
| authority | 1 | career | 3 | conceded | 7 |
| available | 1 | cash | 3 | concentrate | 8 |
| awaited | 4 | cast | 6 | concert | 10 |
| award | 4 | casualties | 10 | concluded | 6 |
| aware | 2 | category | 6 | conclusion | 8 |
| bail | 10 | celebrated | 6 | condemned | 9 |
| ban | 3 | cellphone | 10 | conducted | 3 |
| beach | 4 | challenge | 1 | conference | 1 |
| behalf | 5 | champion | 5 | confirmed | 3 |
| benefits | 3 | chancellor | 10 | conflict | 8 |
| bet | 9 | channel | 7 | confrontation | 7 |
| bid | 1 | chaos | 7 | congress | 10 |
| bomb | 8 | chase | 8 | consecutive | 10 |
| bond | 9 | chip | 8 | consent | 9 |
| bonus | 9 | circuit | 8 | consequences | 3 |
| boom | 4 | circumstances | 9 | conservative | 6 |

| | | | | | |
|---|---|---|---|---|---|
| considerable | 6 | decade | 1 | editor | 5 |
| consistent | 4 | declined | 2 | element | 8 |
| construction | 2 | defence | 1 | eliminate | 8 |
| consultant | 3 | defendants | 9 | email | 10 |
| consumers | 7 | deficit | 5 | embarrassed | 6 |
| contact | 7 | definitely | 4 | embrace | 5 |
| contend | 6 | definition | 5 | emerged | 2 |
| contest | 9 | deliberately | 6 | emotional | 4 |
| contract | 1 | demonstrations | 4 | emphasis | 5 |
| contrast | 6 | denied | 1 | enable | 5 |
| contributed | 2 | departure | 8 | endorsed | 9 |
| controversy | 2 | depressed | 8 | energy | 5 |
| convention | 8 | deputy | 5 | enforcement | 4 |
| converted | 10 | designed | 2 | enormous | 8 |
| convicted | 8 | desperate | 8 | enrichment | 9 |
| convinced | 6 | despite | 1 | enrolled | 10 |
| cooperation | 9 | detention | 10 | ensure | 4 |
| cope | 10 | disabled | 9 | entitled | 8 |
| corporate | 3 | disaster | 5 | environment | 3 |
| counter | 8 | disclose | 7 | equipment | 5 |
| counterparts | 8 | discount | 7 | era | 8 |
| county | 10 | discrimination | 10 | errors | 6 |
| couple | 2 | display | 6 | erupted | 9 |
| create | 1 | dispute | 2 | established | 2 |
| credit | 2 | distinctive | 7 | estimated | 2 |
| crew | 6 | distribution | 6 | evaluate | 9 |
| crisis | 3 | dividend | 10 | eventually | 2 |
| criticised | 2 | document | 5 | evidence | 1 |
| criticism | 6 | domestic | 4 | ex- | 5 |
| crucial | 3 | dominated | 4 | exceed | 8 |
| culture | 4 | donations | 10 | exclude | 9 |
| curb | 9 | draft | 7 | executed | 6 |
| deadline | 4 | dramatic | 6 | executive | 1 |
| debate | 2 | drug | 3 | expand | 2 |
| debut | 9 | echoed | 7 | experts | 4 |

| | | | | | |
|---|---|---|---|---|---|
| exports | 7 | guidelines | 9 | instance | 6 |
| exposed | 4 | guys | 7 | institute | 1 |
| extract | 10 | hail | 9 | intelligence | 7 |
| facilities | 4 | halt | 4 | intense | 2 |
| factor | 4 | haul | 8 | interim | 10 |
| feature | 4 | headlines | 7 | internal | 4 |
| federal | 5 | headquarters | 3 | internet | 3 |
| federation | 4 | height | 8 | interview | 3 |
| feeding | 3 | highlighted | 7 | investigation | 1 |
| fees | 7 | huge | 1 | investment | 1 |
| final | 1 | identified | 3 | involved | 1 |
| financial | 1 | ignore | 4 | isolated | 6 |
| fines | 6 | image | 2 | issue | 1 |
| fled | 9 | immigration | 9 | items | 9 |
| flexible | 8 | impact | 2 | jail | 5 |
| flu | 9 | implications | 7 | job | 1 |
| focus | 1 | import | 5 | journal | 10 |
| forecast | 7 | imposed | 2 | journalists | 2 |
| founder | 5 | impression | 4 | junior | 9 |
| franchise | 10 | incident | 5 | jury | 10 |
| frustrated | 3 | income | 7 | justify | 5 |
| fuel | 1 | incorporated | 8 | keen | 10 |
| fund | 3 | indicated | 2 | kids | 6 |
| fundamental | 5 | individual | 2 | label | 8 |
| generation | 1 | inevitable | 7 | labour | 5 |
| giant | 2 | infection | 9 | lane | 8 |
| global | 5 | inflation | 8 | lap | 9 |
| goal | 3 | infrastructure | 9 | launched | 3 |
| golf | 10 | initial | 2 | league | 7 |
| goods | 6 | initiative | 6 | leaking | 9 |
| grab | 8 | injury | 3 | legal | 3 |
| grade | 8 | insisted | 3 | licence | 8 |
| graduate | 10 | inspector | 6 | link | 5 |
| grant | 2 | inspired | 9 | lobby | 9 |
| guarantee | 4 | installed | 10 | location | 4 |

| | | | | | |
|---|---|---|---|---|---|
| Ltd. | 10 | opponent | 4 | premier | 7 |
| magazine | 7 | optimistic | 5 | previous | 1 |
| maintaining | 2 | option | 1 | primary | 6 |
| major | 1 | outcome | 3 | prime | 1 |
| margin | 5 | overall | 2 | prince | 10 |
| massive | 6 | overnight | 6 | principal | 7 |
| maximum | 3 | pace | 4 | principle | 8 |
| media | 3 | panel | 6 | priority | 7 |
| medical | 2 | panic | 5 | pro- | 8 |
| mental | 7 | participation | 6 | proceedings | 2 |
| military | 3 | partner | 1 | process | 1 |
| minimum | 7 | passion | 7 | professional | 2 |
| ministry | 8 | patients | 10 | profile | 4 |
| minor | 2 | peak | 9 | project | 1 |
| mirror | 7 | penalty | 3 | prominent | 9 |
| mission | 7 | pension | 9 | promote | 4 |
| mobile | 6 | percent | 1 | prop | 10 |
| monitors | 2 | period | 1 | prosecution | 5 |
| mood | 8 | personality | 10 | prospects | 2 |
| motivated | 4 | personnel | 6 | protests | 3 |
| mount | 4 | physical | 6 | province | 7 |
| mutual | 8 | pit | 10 | provoke | 8 |
| negative | 4 | pledged | 7 | publisher | 5 |
| negotiations | 1 | plunged | 9 | purchase | 10 |
| network | 3 | plus | 2 | pursue | 5 |
| nomination | 9 | PM (Prime Minister) | 8 | quit | 7 |
| normal | 2 | policy | 3 | quoted | 7 |
| obligations | 8 | polls | 8 | raid | 7 |
| obstacle | 8 | port | 3 | rally | 9 |
| obviously | 4 | posed | 7 | range | 1 |
| occupation | 9 | positive | 3 | reaction | 3 |
| occupied | 4 | potential | 1 | recalled | 7 |
| occurred | 6 | pre- | 5 | recovery | 2 |
| odd | 6 | predicted | 2 | recruiting | 2 |
| olympic | 5 | | | refinery | 10 |

| | | | | | |
|---|---|---|---|---|---|
| regain | 8 | scheduled | 4 | structure | 2 |
| regime | 7 | scored | 1 | stunned | 5 |
| region | 3 | section | 4 | style | 4 |
| register | 6 | security | 1 | suburb | 9 |
| regulation | 7 | seeking | 1 | successor | 10 |
| rejected | 3 | select | 2 | sufficient | 9 |
| release | 1 | senior | 1 | sum | 9 |
| reluctant | 5 | series | 1 | super | 7 |
| relying | 2 | session | 4 | supreme | 9 |
| remote | 9 | sex | 8 | surge | 8 |
| removed | 4 | shareholders | 9 | surgery | 7 |
| required | 1 | shift | 4 | survive | 2 |
| research | 5 | significant | 1 | suspended | 6 |
| residents | 5 | similar | 1 | sustained | 6 |
| resolve | 9 | site | 1 | switch | 6 |
| resort | 9 | slim | 10 | symbol | 6 |
| resource | 5 | slumped | 10 | tackle | 5 |
| respond | 2 | smart | 6 | tactics | 6 |
| response | 4 | soared | 7 | tank | 9 |
| restrictions | 4 | sole | 5 | tape | 8 |
| resume | 9 | source | 2 | target | 3 |
| retain | 3 | sparked | 10 | task | 7 |
| revealed | 2 | specific | 4 | team | 1 |
| revenue | 5 | speculation | 4 | technical | 4 |
| reverse | 2 | spokesman | 1 | technology | 3 |
| revised | 8 | spokeswoman | 6 | teenagers | 6 |
| revolution | 9 | sponsored | 4 | tensions | 4 |
| riot | 8 | spurs | 9 | territory | 7 |
| role | 1 | stability | 3 | terrorist | 5 |
| route | 4 | stake | 5 | testified | 9 |
| routine | 3 | statistics | 4 | testimony | 8 |
| rumors | 6 | status | 4 | text | 10 |
| sanctions | 7 | strain | 6 | theme | 8 |
| scandal | 7 | strategy | 1 | tiny | 7 |
| scared | 6 | stress | 4 | toll | 9 |

| traditional | 2 |
|---|---|
| traffic | 8 |
| transfer | 2 |
| transformed | 5 |
| transmission | 9 |
| transport | 5 |
| trend | 7 |
| tribunal | 10 |
| triple | 10 |
| ultimately | 2 |
| undermine | 6 |
| unique | 8 |
| urgent | 8 |
| utility | 8 |
| variable | 5 |
| vast | 5 |
| venture | 6 |
| verdict | 10 |
| version | 4 |
| veteran | 6 |
| veto | 10 |
| vice | 6 |
| victim | 5 |
| video | 5 |
| virus | 9 |
| vital | 4 |
| volume | 6 |
| volunteers | 10 |
| vulnerable | 8 |
| watchdog | 8 |
| website | 5 |
| widespread | 5 |
| withdrawal | 2 |
| zone | 2 |