

Learning Vocabulary in Another Language

I. S. P. Nation

Victoria University of Wellington



CAMBRIDGE
UNIVERSITY PRESS

PUBLISHED BY THE PRESS SYNDICATE OF THE UNIVERSITY OF CAMBRIDGE
The Pitt Building, Trumpington Street, Cambridge, United Kingdom

CAMBRIDGE UNIVERSITY PRESS
The Edinburgh Building, Cambridge CB2 2RU, UK
40 West 20th Street, New York, NY 10011-4211, USA
10 Stamford Road, Oakleigh, VIC 3166, Australia
Ruiz de Alarcón 13, 28014 Madrid, Spain
Dock House, The Waterfront, Cape Town 8001, South Africa

<http://www.cambridge.org>

© Cambridge University Press 2001

This book is in copyright. Subject to statutory exception
and to the provisions of relevant collective licensing agreements,
no reproduction of any part may take place without
the written permission of Cambridge University Press.

First published 2001

Printed in the United Kingdom at the University Press, Cambridge

Typeface Sabon 10.5/12 pt. System QuarkXPress™ [SE]

A catalogue record for this book is available from the British Library

ISBN 0 521 800927 hardback
ISBN 0 521 804981 paperback

Contents

Series editors' preface	<i>page</i> xiii
Acknowledgements	xiv
Introduction	1
Learning goals	1
The four strands	2
Main themes	3
The audience for this book	4
1 The goals of vocabulary learning	6
How much vocabulary do learners need to know?	6
How many words are there in the language?	6
How many words do native speakers know?	8
How much vocabulary do you need to use another language?	9
High-frequency words	13
Specialised vocabulary	17
Low-frequency words	19
Testing vocabulary knowledge	21
2 Knowing a word	23
Learning burden	23
The receptive/productive distinction	24
The scope of the receptive/productive distinction	26
Experimental comparisons of receptive and productive vocabulary	30
Aspects of knowing a word	33
Levelt's process model of language use	34
Spoken form	40
Written form	44
Word parts	46
Connecting form and meaning	47
Concept and referents	49
	vii

	Associations	52
	Grammatical functions	55
	Collocations	56
	Constraints on use	57
	Item knowledge and system knowledge	58
3	Teaching and explaining vocabulary	60
	Learning from teaching and learning activities	60
	Vocabulary in classrooms	74
	Repetition and learning	74
	Communicating meaning	81
	Helping learners comprehend and learn from definitions	90
	Spending time on words	93
	Rich instruction	94
	Arguments against rich instruction	95
	Providing rich instruction	97
	Spoken form	98
	Written form	98
	Word parts	100
	Strengthening the form–meaning connection	101
	Concept and referents	102
	Associations	104
	Grammar	106
	Collocation	106
	Constraints on use	106
	Vocabulary teaching procedures	107
	Computer-assisted vocabulary learning	108
	Using concordances	111
	Research on CAVL	112
4	Vocabulary and listening and speaking	114
	What vocabulary knowledge is needed for listening?	114
	Providing vocabulary support for listening	116
	Learning vocabulary from listening to stories	117
	Learning vocabulary through negotiation	123
	The vocabulary of speaking	125
	Developing fluency with spoken vocabulary	127
	Using teacher input to increase vocabulary knowledge	129
	Using labelled diagrams	131
	Using cooperative tasks to focus on vocabulary	133
	How can a teacher design activities to help incidental vocabulary learning?	134
	Designing and adapting activities	139

5	Vocabulary and reading and writing	144
	Vocabulary size and successful reading	144
	Learning vocabulary through reading	149
	Vocabulary and extensive reading	150
	Extensive reading by non-native speakers of texts written for young native speakers	151
	Extensive reading with graded readers	154
	Extensive reading of unsimplified texts	154
	Extensive reading and vocabulary growth	155
	Intensive reading and direct teaching	156
	Preteaching	157
	Vocabulary exercises with reading texts	158
	Analysis of vocabulary exercises	159
	Readability	161
	What are graded readers?	162
	Designing and using a simplified reading scheme for vocabulary development	164
	How to simplify	171
	Alternatives to simplification	173
	Glossing	174
	Vocabulary and the quality of writing	177
	Measures of vocabulary size and growth in writing	178
	Bringing vocabulary into productive use	180
	Responding to vocabulary use in written work	185
6	Specialised uses of vocabulary	187
	Academic vocabulary	187
	The importance of academic vocabulary	189
	Making an academic vocabulary list	191
	Sequencing the introduction of academic vocabulary	193
	The nature and role of academic vocabulary	194
	Testing academic vocabulary	196
	Learning academic vocabulary	196
	Technical vocabulary	198
	Distinguishing technical vocabulary from other vocabulary	198
	Making lists of technical vocabulary	201
	Learning technical vocabulary	203
	Vocabulary in discourse	205
	Vocabulary and the information content of the text	206
	Vocabulary and the relationship between the writer or speaker and reader or listener	209
	Vocabulary and the organisation of the text	210
	Words in discourse	214

7	Vocabulary learning strategies and guessing from context	217
	A taxonomy of vocabulary learning strategies	217
	Planning vocabulary learning	218
	Sources: finding information about words	219
	Processes: establishing vocabulary knowledge	221
	Training in strategy choice and use	222
	Learners' use of strategies	224
	Procedures that integrate strategies	229
	Learning words from context	232
	Intentional and incidental learning	232
	What proportion of unknown words can be guessed from context?	233
	How much vocabulary is learned from context?	236
	What can be learned from context?	240
	What clues does a context provide and how effective are they?	242
	What are the causes of poor guessing?	246
	Do different learners approach guessing in the same way?	247
	How can teachers help learners improve learning from context?	250
	How can learners be trained to guess from context?	250
	Learning from context and attention-drawing activities	251
	Do glossing and dictionary use help vocabulary learning?	252
	Formats for testing or practising guessing	253
	Steps in the guessing-from-context strategy	256
	Training learners in the strategy of guessing from context	261
8	Word study strategies	263
	Word parts	263
	Is it worthwhile learning word parts?	264
	Studies of the sources of English vocabulary	264
	Studies of the proportion of affixed words	265
	Studies of the frequency of affixes	266
	Do language users see words as being made of parts?	269
	Word stems	272
	The knowledge required to use word parts	272
	Monitoring and testing word building skills	275
	The word part strategy	278
	Using dictionaries	281
	Is it necessary or worth training learners to use dictionaries?	282
	What skills are needed to use a dictionary?	284

What dictionaries are the best?	288
Evaluating dictionaries	290
Dictionary use and learning	296
Learning from word cards	296
Criticisms of direct vocabulary learning	297
Decontextualised learning and memory	297
Decontextualised learning and use	299
The contribution of decontextualised learning	301
The values of learning from word cards	302
The word card strategy	303
Training learners in the use of word cards	315
9 Chunking and collocation	317
Chunking	319
The advantages and disadvantages of chunking	320
Language knowledge is collocational knowledge	321
Fluent and appropriate language use requires collocational knowledge	323
Some words occur in a limited set of collocations	324
Classifying collocations	328
The evidence for collocation	333
Collocation and teaching	335
Encouraging chunking	336
Chunking through fluency development	336
Chunking through language-focused attention	340
Memorising unanalysed chunks	343
10 Testing vocabulary knowledge and use	344
What kind of vocabulary test is the best?	344
Is it enough to ask learners if they know the word?	346
Should choices be given?	349
Should translations be used?	351
Should words be tested in context?	352
How can depth of knowledge about a word be tested?	354
How can we measure words that learners don't know well?	358
How can we measure how well learners actually use words?	361
How can we measure total vocabulary size?	362
Choosing a test item type	372
Types of tests	373
How can we test to see where learners need help?	373
How can we test whether a small group of words in a course has been learned?	374

How can we test whether the total vocabulary of the course has been learned?	375
How can we measure how well learners have control of the important vocabulary learning strategies?	378
11 Designing the vocabulary component of a language course	380
Goals	380
Needs analysis	381
Environment analysis	383
Principles of vocabulary teaching	384
Content and sequencing	385
Format and presentation	388
Monitoring and assessment	389
Evaluation	391
Autonomy and vocabulary learning	394
Appendixes	
1. Headwords of the <i>Academic Word List</i>	407
2. 1,000 word level tests	412
3. A Vocabulary Levels Test: Test B	416
4. Productive Levels Test: Version C	425
5. Vocabulary Levels Dictation Test	429
6. Function words	430
References	432
Subject index	464
Author index	470

1 *The goals of vocabulary learning*

How much vocabulary do learners need to know?

Whether designing a language course or planning our own course of study, it is useful to be able to set learning goals that will allow us to use the language in the ways we want to. When we plan the vocabulary goals of a long-term course of study, we can look at three kinds of information to help decide how much vocabulary needs to be learned: the number of words in the language, the number of words known by native speakers and the number of words needed to use the language.

How many words are there in the language?

The most ambitious goal is to know all of the language. However, even native speakers do not know all the vocabulary of the language. There are numerous specialist vocabularies, such as those of nuclear physics or computational linguistics, which are known only by the small groups who specialise in those areas. Still, it is interesting to have some idea of how many words there are in the language. This is not an easy question to resolve because there are numerous other questions which affect the way we answer it, including the following.

What do we count as a word? Do we count *book* and *books* as the same word? Do we count *green* (the colour) and *green* (a large grassed area) as the same word? Do we count people's names? Do we count the names of products like *Fab*, *Pepsi*, *Vegemite*, *Chevrolet*? The few brave or foolish attempts to answer these questions and the major question 'How many words are there in English?' have counted the number of words in very large dictionaries. *Webster's Third New International Dictionary* is the largest non-historical dictionary of English. It contains around 114,000 word families excluding proper names (Goulden, Nation and Read, 1990). This is a very large number and is well beyond the goals of most first and second language learners.

There are several ways of deciding what words will be counted.

Tokens

One way is simply to count every word form in a spoken or written text and if the same word form occurs more than once, then each occurrence of it is counted. So the sentence 'It is not easy to say it correctly' would contain eight words, even though two of them are the same word form, *it*. Words which are counted in this way are called 'tokens', and sometimes 'running words'. If we try to answer questions like 'How many words are there on a page or in a line?' 'How long is this book?' 'How fast can you read?' 'How many words does the average person speak per minute?' then our unit of counting will be the token.

Types

We can count the words in the sentence 'It is not easy to say it correctly' another way. If we see the same word again, we do not count it again. So the sentence of eight tokens consists of seven different words or 'types'. We count words in this way if we want to answer questions like 'How large was Shakespeare's vocabulary?' 'How many words do you need to know to read this book?' 'How many words does this dictionary contain?'

Lemmas

A lemma consists of a headword and some of its inflected and reduced (*n't*) forms. Usually, all the items included under a lemma are the same part of speech (Francis and Kučera, 1982: 461). The English inflections consist of plural, third person singular present tense, past tense, past participle, *-ing*, comparative, superlative and possessive (Bauer and Nation, 1993). The Thorndike and Lorge (1944) frequency count used lemmas as the basis for counting, and the more recent computerised count on the *Brown Corpus* (Francis and Kučera, 1982) has produced a lemmatised list. In the Brown count the comparative and superlative forms are not included in the lemma, and the same form used as a different part of speech (*walk* as a noun, *walk* as a verb) are not in the same lemma. Variant spellings (*favor*, *favour*) are usually included as part of the same lemma when they are the same part of speech.

Lying behind the use of lemmas as the unit of counting is the idea of learning burden (Swenson and West, 1934). The learning burden of an item is the amount of effort required to learn it. Once learners can use the inflectional system, the learning burden of for example *mends*, if

the learner already knows *mend*, is negligible. One problem in forming lemmas is to decide what will be done with irregular forms such as *mice*, *is*, *brought*, *beaten* and *best*. The learning burden of these is clearly heavier than the learning burden of regular forms like *books*, *runs*, *talked*, *washed* and *fastest*. Should the irregular forms be counted as a part of the same lemma as their base word or should they be put into separate lemmas? Lemmas also separate closely related items such as the adjective and noun uses of words like *original*, and the noun and verb uses of words like *display*. An additional problem with lemmas is what is the headword – the base form or the most frequent form? (Sinclair, 1991: 41-42).

Using the lemma as the unit of counting greatly reduces the number of units in a corpus. Bauer and Nation (1993) calculate that the 61,805 tagged types (or 45,957 untagged types) in the *Brown Corpus* become 37,617 lemmas which is a reduction of almost 40% (or 18% for untagged types). Nagy and Anderson (1984) estimated that 19,105 of the 86,741 types in the Carroll, Davies and Richman (1971) corpus were regular inflections.

Word families

Lemmas are a step in the right direction when trying to represent learning burden in the counting of words. However, there are clearly other affixes which are used systematically and which greatly reduce the learning burden of derived words containing known base forms. These include affixes like *-ly*, *-ness* and *un-*. A word family consists of a headword, its inflected forms, and its closely related derived forms.

The major problem in counting using word families as the unit is to decide what should be included in a word family and what should not. Learners' knowledge of the prefixes and suffixes develops as they gain more experience of the language. What might be a sensible word family for one learner may be beyond another learner's present level of proficiency. This means that it is usually necessary to set up a scale of word families, starting with the most elementary and transparent members and moving on to less obvious possibilities.

How many words do native speakers know?

A less ambitious way of setting vocabulary learning goals is to look at what native speakers of the language know. Unfortunately, research on measuring vocabulary size has generally been poorly done (Nation, 1993c), and the results of the studies stretching back to the late nine-

teenth century are often wildly incorrect. We will look at the reasons for this later in this book.

Recent reliable studies (Goulden, Nation and Read, 1990; Zechmeister, Chronis, Cull, D'Anna and Healy, 1995) suggest that educated native speakers of English know around 20,000 word families. These estimates are rather low because the counting unit is word families which have several derived family members and proper nouns are not included in the count. A very rough rule of thumb would be that for each year of their early life, native speakers add on average 1,000 word families a year to their vocabulary. These goals are manageable for non-native speakers of English, especially those learning English as a second rather than foreign language, but they are way beyond what most learners of English as another language can realistically hope to achieve.

How much vocabulary do you need to use another language?

Studies of native speakers' vocabulary seem to suggest that second language learners need to know very large numbers of words. While this may be useful in the long term, it is not an essential short-term goal. This is because studies of native speakers' vocabulary growth see all words as being of equal value to the learner. Frequency based studies show very strikingly that this is not so, and that some words are much more useful than others.

Table 1.1 shows part of the results of a frequency count of just under 500 running words of the Ladybird version of *The Three Little Pigs*. It contains 124 different word types.

Note the large proportion of words occurring only once, and the very high frequency of the few most frequent words. When we look at texts our learners may have to read and conversations that are like ones that they may be involved in, we find that a relatively small amount of well-chosen vocabulary can allow learners to do a lot. To see this, let us look at an academic reading text and examine the different kinds of vocabulary it contains. The text is from Neville Peat's (1987) *Forever the Forest. A West Coast Story* (Hodder and Stoughton, Auckland).

Sustained-yield management ought to be long-term government **policy** in *indigenous* forests *zoned* for production. The adoption of such a **policy** would represent a *breakthrough* – the boundary between a *pioneering*, **extractive phase** and an *era* in which the *timber* industry **adjusted** to living with the forests in *perpetuity*. A forest **sustained** is a forest in which harvesting and *mortality* combined do not **exceed** *regeneration*. Naturally enough, faster-growing forests produce more *timber*, which is why attention

Table 1.1. *An example of the results of a frequency count*

the	41	met	3	come	1
little	25	myself	3	door	1
pig	22	not	3	down	1
house	17	on	3	fell	1
a	16	pigs	3	go	1
and	16	please	3	grew	1
said	14	pleased	3	had	1
he	12	shall	3	hair	1
I	10	soon	3	here	1
me	10	stronger	3	him	1
some	9	that	3	houses	1
wolf	9	they	3	huff	1
build	8	three	3	knocked	1
't	8	want	3	live	1
third	8	who	3	long	1
was	8	with	3	mother	1
of	7	won	3	must	1
straw	7	yes	3	my	1
to	7	yours	3	next	1
you	7	big	2	off	1
man	6	by	2	once	1
second	6	care	2	one	1
catch	5	chin	2	puff	1
first	5	day	2	road	1
for	5	does	2	set	1
will	5	huffed	2	so	1
bricks	4	let	2	their	1
built	4	'm	2	them	1
himself	4	no	2	there	1
now	4	puffed	2	took	1
sticks	4	strong	2	up	1
than	4	take	2	upon	1
very	4	then	2	us	1
asked	3	time	2	walked	1
carrying	3	too	2	we	1
eat	3	along	1	went	1
gave	3	are	1	were	1
give	3	ate	1	which	1
his	3	blow	1	your	1
in	3	but	1	yourselves	1
it	3	came	1		
'll	3	chinny	1		

would tend to swing from *podocarps* to *beech* forests regardless of the state of the *podocarp resource*. The colonists cannot be blamed for *plunging* in without thought to whether the **resource** had limits. They brought from *Britain* little experience or understanding of how to **maintain** forest **structure** and a *timber* supply for all time. Under *German* management it might have been different here. The *Germans* have practised the **sustained approach** since the seventeenth century when they faced a *timber* shortage as a result of a **series** of wars. In *New Zealand* in the latter part of the twentieth century, an **anticipated** shortage of the most valuable native *timber, rimu*, prompts a **similar response** – no more **contraction** of the *indigenous* forest and a balancing of yield with *increment* in **selected areas**.

This is not to say the idea is being *aired* here for the first time. Over a century ago the first *Conservator* of Forests proposed **sustained** harvesting. He was cried down. There were far too many trees left to bother about it. And yet in the *pastoral context* the dangers of *overgrazing* were **appreciated** early in the piece. *New Zealand geography* students are taught to this day how *overgrazing* causes the *degradation* of the soil and hillsides to slide away, and that with them can go the *viability* of hill-country sheep and cattle farming. That a forest could be *overgrazed* as easily was not widely accepted until much later – so late, in fact, that the *counter* to it, **sustained-yield** management, would be forced upon the industry and come as a shock to it. It is a simple enough **concept** on paper: balance harvest with growth and you have a natural *renewable resource*; forest products forever. **Plus** the social and **economic benefits** of regular work and **income**, a regular *timber* supply and relatively **stable** markets. **Plus** the **environmental benefits** that *accrue* from **minimising** the **impact** on soil and water qualities and wildlife.

In practice, however, **sustainability** depends on how well the **dynamics** of the forest are understood. And these **vary** from **area** to **area** according to forest make-up, soil *profile*, *altitude*, *climate* and **factors** which forest science may yet discover. *Ecology* is deep-felt.

We can distinguish four kinds of vocabulary in the text: high-frequency words (unmarked in the text), academic words (in bold), and technical and low-frequency words (in italics).

High-frequency words

In the example text, these words are not marked at all and include function words: *in, for, the, of, a*, etc. Appendix 6 contains a complete list of function words. The high-frequency words also include many content words: *government, forests, production, adoption, represent, boundary*. The classic list of high-frequency words is Michael West's (1953a) *A General Service List of English Words* which contains around 2,000 word families. Almost 80% of the running words in the text are high-frequency words.

Academic words

The text is from an academic textbook and contains many words that are common in different kinds of academic texts: *policy, phase, adjusted, sustained*. Typically these words make up about 9% of the running words in the text. The best list of these is the *Academic Word List* (Coxhead, 1998). Appendix 1 contains the 570 headwords of this list. This small list of words is very important for anyone using English for academic purposes (see chapter 6).

Technical words

The text contains some words that are very closely related to the topic and subject area of the text. These words include *indigenous, regeneration, podocarp, beech, rimu* (a New Zealand tree) and *timber*. These words are reasonably common in this topic area but not so common elsewhere. As soon as we see them we know what topic is being dealt with. Technical words like these typically cover about 5% of the running words in a text. They differ from subject area to subject area. If we look at technical dictionaries, such as dictionaries of economics, geography or electronics, we usually find about 1,000 entries in each dictionary.

Low-frequency words

The fourth group is the low-frequency words. Here, this group includes words like *zoned, pioneering, perpetuity, aired* and *pastoral*. They make up over 5% of the words in an academic text. There are thousands of them in the language, by far the biggest group of words. They include all the words that are not high-frequency words, not academic words and not technical words for a particular subject. They consist of technical words for other subject areas, proper nouns, words that almost got into the high-frequency list, and words that we rarely meet in our use of the language.

Let us now look at a longer text and a large collection of texts.

Sutarsyah, Nation and Kennedy (1994) looked at a single economics textbook to see what vocabulary would be needed to read the text. The textbook was 295,294 words long. Table 1.2 shows the results. The academic word list used in the study was the *University Word List* (Xue and Nation, 1984).

What should be clear from this example and from the text looked at earlier is that a reasonably small number of words covers a lot of text.

Table 1.2. *Text coverage by the different kinds of vocabulary in an economics textbook*

Type of vocabulary	Number of words	Text coverage
1st 2000 word families	1,577	82.5%
Academic vocabulary	636	8.7%
Other vocabulary	3,225	8.8%
Total	5,438	100.0%

Table 1.3. *The coverage by the different kinds of vocabulary in an academic corpus*

Type of vocabulary	% coverage
1st 1000 words	71.4%
2nd 1000 words	4.7%
Academic Word List (570 words)	10.0%
Others	13.9%
Total	100.0%

Coxhead (1998) used an academic corpus made up of a balance of science, arts, commerce and law texts totalling 3,500,000 running words. Table 1.3 gives the coverage figures for this corpus.

Figure 1.1 presents the proportions in a diagrammatic form. The size of each of the sections of the right-hand box indicates the proportion of the text taken up by each type of vocabulary.

Table 1.4 gives the typical figures for a collection of texts consisting of five million running words.

Some very important generalisations can be drawn from Table 1.4 and the other information that we have looked at. We will look at these generalisations and at questions that they raise. Brief answers to the questions will be given here but will be examined much more closely in later chapters.

High-frequency words

There is a small group of high-frequency words which are very important because these words cover a very large proportion of the running words in spoken and written texts and occur in all kinds of uses of the language.

Sustained-yield management ought to be long-term government **policy** in *indigenous* forests **zoned** for production. The adoption of such a **policy** would represent a *breakthrough* – the boundary between a *pioneering, extractive phase* and an *era* in which the *timber* industry **adjusted** to living with the forests in *perpetuity*. A forest **sustained** is a forest in which harvesting and *mortality* combined do not **exceed** *regeneration*. Naturally enough, faster-growing forests produce more *timber*, which is why attention would tend to swing from *podocarps* to *beech* forests regardless of the state of the *podocarp* **resource**. The colonists cannot be blamed for *plunging* in without thought to whether the **resource** had limits. They brought from *Britain* little experience or understanding of how to **maintain** forest **structure** and a *timber* supply for all time. Under *German* management it might have been different here. The *Germans* have practised the **sustained approach** since the seventeenth century when they faced a *timber* shortage as a result of a *series* of wars. In *New Zealand* in the latter part of the twentieth century, an **anticipated** shortage of the most valuable native *timber, rimu*, prompts a **similar response** – no more **contraction** of the *indigenous* forest and a balancing of yield with **increment** in **selected areas**.

This is not to say the idea is being *aired* here for the first time. Over a century ago the first *Conservator* of Forests proposed **sustained** harvesting. He was *cried down*. There were far too many trees left to bother about it. And yet in the *pastoral context* the dangers of *overgrazing* were **appreciated** early in the piece. *New Zealand geography* students are taught to this day how *overgrazing* causes the *degradation* of the soil and hillsides to slide away, and that with them can go the *viability* of hill-country sheep and cattle farming. That a forest could be *overgrazed* as easily was not widely accepted until much later – so late, in fact, that the *counter* to it, **sustained-yield** management, would be forced upon the industry and come as a shock to it.

<p>High-frequency vocabulary</p> <p>2000 words 80% or more text coverage a, equal, places, <i>behaves</i>, <i>educate</i></p>
<p>Academic vocabulary</p>
<p>Technical vocabulary</p>
<p>Low-frequency vocabulary</p>

Figure 1.1 Vocabulary type and coverage in an academic text

How large is this group of words? The usual way of deciding how many words should be considered as high-frequency words is to look at the text coverage provided by successive frequency-ranked groups of words. The teacher or course designer then has to decide where the coverage gained by spending teaching time on these words is no longer worthwhile. Table 1.5 shows coverage figures for each successive 1,000 lemmas from the *Brown Corpus* – a collection of various 2,000-word texts of American English totalling just over one million tokens.

Usually the 2,000-word level has been set as the most suitable limit for high-frequency words. Nation and Hwang (1995) present

Table 1.4. *Vocabulary size and coverage (Carroll, Davies and Richman (1971))*

Number of words	% text coverage
86,741	100
43,831	99
12,448	95
5,000	89.4
4,000	87.6
3,000	85.2
2,000	81.3
1,000	74.1
100	49
10	23.7

Table 1.5. *The percentage text coverage of each successive 1000 lemmas in the Brown Corpus*

1000 word (lemma) level	% coverage of text (tokens)
1000	72
2000	79.7
3000	84
4000	86.7
5000	88.6
6000	89.9

evidence that counting the 2,000 most frequent words of English as the high-frequency words is still the best decision for learners going on to academic study.

What are the words in this group? As has been noted, the classic list of high-frequency words is Michael West's *General Service List* which contains 2,000 word families. About 165 word families in this list are function words such as *a, some, two, because* and *to* (see appendix 6). The rest are content words, that is nouns, verbs, adjectives and adverbs. Older series of graded readers are based on this list.

How stable are the high-frequency words? In other words, does one properly researched list of high-frequency words differ greatly from another? Frequency lists may disagree with each other about the frequency rank order of particular words but if the research is based on a well-designed corpus there is generally about 80% agreement about

Table 1.6. *Ways of learning and teaching high-frequency words*

Direct teaching	Teacher explanation Peer teaching
Direct learning	Study from word cards Dictionary use
Incidental learning	Guessing from context in extensive reading Use in communication activities
Planned encounters	Graded reading Vocabulary exercises

what particular words should be included. Nation and Hwang's (1995) research on the *General Service List* showed quite large overlap between it and more recent frequency counts. Replacing some of the words in the *General Service List* with other words resulted in only a 1% increase in coverage. It is important to remember that the 2,000 high-frequency words of English consist of some words that have *very* high frequencies and some words that are only slightly more frequent than others not in the list. The first 1,000 words cover about 77% and the second 1,000 about 5% of the running words in academic texts. When making a list of high-frequency words, both frequency and range must be considered. Range is measured by seeing how many different texts or subcorpora each particular word occurs in. A word with wide range occurs in many different texts or subcorpora.

How should teachers and learners deal with these words? The high-frequency words of the language are clearly so important that considerable time should be spent on them by teachers and learners. The words are a small enough group to enable most of them to get attention over the span of a long-term English programme. This attention can be in the form of direct teaching, direct learning, incidental learning, and planned meetings with the words. The time spent on them is well justified by their frequency, coverage and range. Table 1.6 lists some of the teaching and learning possibilities that will be explored in more detail in later chapters of this book.

In general, high-frequency words are so important that anything that teachers and learners can do to make sure they are learned is worth doing.

Table 1.7. *Text type and text coverage by the most frequent 2000 words of English and an academic word list in four different kinds of texts*

Levels	Conversation	Fiction	Newspapers	Academic text
1st 1000	84.3%	82.3%	75.6%	73.5%
2nd 1000	6%	5.1%	4.7%	4.6%
Academic	1.9%	1.7%	3.9%	8.5%
Other	7.8%	10.9%	15.7%	13.3%

Specialised vocabulary

It is possible to make specialised vocabularies which provide good coverage for certain kinds of texts. These are a way of extending the high-frequency words for special purposes.

What special vocabularies are there? Special vocabularies are made by systematically restricting the range of topics or language uses investigated. It is thus possible to have special vocabularies for speaking, for reading academic texts, for reading newspapers, for reading children's stories, or for letter writing. Technical vocabularies are also specialised vocabularies. Some specialised vocabularies are made by doing frequency counts using a specialised corpus, others are made by experts in the field gathering what they consider to be relevant vocabulary.

There is a very important specialised vocabulary for second language learners intending to do academic study in English. This is the *Academic Word List* (Coxhead, 1998; see appendix 1). It consists of 570 word families that are not in the most frequent 2,000 words of English but which occur reasonably frequently over a very wide range of academic texts; the list is not restricted to a specific discipline. That means that the words are useful for learners studying humanities, law, science or commerce. Academic vocabulary has sometimes been called sub-technical vocabulary because it does not contain technical words but rather formal vocabulary.

The importance of this vocabulary can be seen in the coverage it provides for various kinds of texts (Table 1.7).

Adding the academic vocabulary from the *UWL* to the high-frequency words changes the coverage of academic text from 78.1% to 86.6%. Expressed another way, with a vocabulary of 2,000 words, approximately one word in every five will be unknown. With a vocabulary of 2,000 words plus the *Academic Word List*, approximately

one word in every ten will be unknown. This is a very significant change. If, instead of learning the vocabulary of the *Academic Word List*, the learner had moved on to the third 1,000 most frequent words, instead of an additional 10% coverage there would only have been 4.3% extra coverage.

What kinds of words do they contain? The *Academic Word List* is in appendix 1. Much research remains to be done on this list to explain why the same group of words frequently occur across a very wide range of academic texts. Sometimes a few of them are closely related to the topic, but most probably occur because they allow academic writers to do the things that academic writers do. That is, they allow writers to refer to others' work (*assume, establish, indicate, conclude, maintain*); and they allow writers to work with data in academic ways (*analyse, assess, concept, definition, establish, categories, seek*). We consider this issue again in chapter 6.

Technical words contain a variety of types which range from words that do not usually occur in other subject areas (*cabotage, amortisation*) to those that are formally like high-frequency words but which have specialised meanings (*demand, supply, cost* as used in economics). Chapter 6 looks more fully at technical words.

How large are they? There has been no survey done of the size of technical vocabularies and little research on finding a consistently applied operational definition of what words are technical words. A rough guess from looking at dictionaries of technical vocabulary, such as those for geography, biology and applied linguistics, is that they each contain less than a thousand words.

How can you make a special vocabulary? The *Academic Word List* was made by deciding on the high-frequency words of English and then examining a range of academic texts to find what words were not among the high-frequency words but had wide range and reasonable frequency of occurrence. Range was important because academic vocabulary is intended for general academic purposes. Making a technical vocabulary is a little more problematic. One of the problem areas is that some technical vocabulary occurs in the high-frequency words and the *Academic Word List*. *Wall* in biology, and *price, cost, demand* in economics are all high-frequency words which have particular technical uses. Sutarsyah, Nation and Kennedy (1994) found that 33 content words made up over 10% of the running words of an economics text, but accounted for less than 1% of the running words in a similar sized set of mixed academic texts. One way of making a technical vocabulary is to compare the frequency of words in a specialised text with their frequency in a general corpus.

What should teachers and learners do about specialised vocabulary? Where possible, specialised vocabulary should be treated like high-frequency vocabulary. That is, it should be taught and studied in a variety of complementary ways. Where technical vocabulary is also high-frequency vocabulary, learners should be helped to see the connections and differences between the high-frequency meanings and the technical uses. For example, what is similar between a cell *wall* and other less specialised uses of *wall*? Where technical vocabulary requires specialist knowledge of the field, teachers should train learners in strategies which will help them understand and remember the words. Much technical vocabulary will only make sense in the context of learning the specialised subject matter. Learning the meaning of the technical term *morpheme* needs to be done as a part of the study of linguistics, not before the linguistics course begins.

Low-frequency words

There is a very large group of words that occur very infrequently and cover only a small proportion of any text.

What kinds of words are they?

1. Some low-frequency words are words of moderate frequency that did not manage to get into the high-frequency list. It is important to remember that the boundary between high-frequency and low-frequency vocabulary is an arbitrary one. Any of several thousand low-frequency words could be candidates for inclusion within the high-frequency list simply because their position on a rank frequency list which takes account of range is dependent on the nature of the corpus the list is based on. A different corpus would lead to a different ranking particularly among words on the boundary. This, however, should not be seen as a justification for large amounts of teaching time being spent on low-frequency words at the third or fourth thousand word level. Here are some words that in the *Brown Corpus* fall just outside the high-frequency boundary: *curious*, *wing*, *arm* (vb), *gate*, *approximately*.
2. Many low-frequency words are proper names. Approximately 4% of the running words in the *Brown Corpus* are words like *Carl*, *Johnson* and *Ohio*. In some texts, such as novels and newspapers, proper nouns are like technical words – they are of high-frequency in particular texts but not in others, their meaning is closely related to the message of the text, and they could not be

sensibly pre-taught because their use in the text reveals their meaning. Before you read a novel, you do not need to learn the characters' names.

3. 'One person's technical vocabulary is another person's low-frequency word.' This ancient vocabulary proverb makes the point that, beyond the high-frequency words of the language, people's vocabulary grows partly as a result of their jobs, interests and specialisations. The technical vocabulary of our personal interests is important to us. To others, however, it is not important and from their point of view is just a collection of low-frequency words.
4. Some low-frequency words are simply low-frequency words. That is, they are words that almost every language user rarely uses, for example: *eponymous*, *gibbous*, *bifurcate*, *plummet*, *ploy*. They may represent a rarely expressed idea; they may be similar in meaning to a much more frequent word or phrase; they may be marked as being old-fashioned, very formal, belonging to a particular dialect, or vulgar, or they may be foreign words.

How many low-frequency words are there and how many do learners need to know? A critical issue in answering this question is to decide what will be counted as a word. For the purpose of providing a brief answer to the question of desirable vocabulary size, word families will be used as the unit of counting. Webster's *Third New International Dictionary* (Gove, 1963) contains 267,000 entries of which 113,161 can be counted as base words including base proper words, base compound words, and homographs with unrelated meanings (Goulden, Nation and Read, 1990: 351). Calculations from *The American Heritage Word Frequency Book* (Carroll, Davies and Richman, 1971) suggest that in printed school English there are 88,533 distinct word families (Nagy and Anderson, 1984: 315). Although not all these words need to be known to be a very successful language user, it is very important that learners continue to increase their vocabulary size. To read with minimal disturbance from unknown vocabulary, language users probably need a vocabulary of 15,000 to 20,000 words.

How should teachers and learners deal with low-frequency vocabulary? Teachers' and learners' aims differ with low-frequency vocabulary. The teacher's aim is to train learners in the use of strategies to deal with such vocabulary. These strategies include guessing from context clues, using word parts to help remember words, using vocabulary cards and dictionaries. When teachers spend time on low-frequency

Table 1.8. *The differing focus of teachers' and learners' attention to high- and low-frequency words*

	High-frequency words	Low-frequency words
Attention to each word	Teacher and learners	Learners
Attention to strategies	Teacher and learners	Teacher and learners

words in class, they should be using the words as an excuse for working on the strategies. The learners' aim is to continue to increase their vocabulary. The strategies provide a means of doing this.

As Table 1.8 shows, learners should begin training in the strategies for dealing with vocabulary while they are learning the high-frequency words of the language. When learners know the high-frequency vocabulary and move to the study of low-frequency words, the teacher does not spend substantial amounts of class time explaining and giving practice with vocabulary, but instead concentrates on expanding and refining the learners' control of vocabulary learning and coping strategies. Learners however should continue to learn new words.

Testing vocabulary knowledge

In this chapter, a very important distinction has been made between high-frequency words and low-frequency words. This distinction has been made on the basis of the frequency, coverage and quantity of these words. The distinction is important because teachers need to deal with these two kinds of words in quite different ways, and teachers and learners need to ensure that the high-frequency words of the language are well known.

It is therefore important that teachers and learners know whether the high-frequency words have been learned. Appendix 3 of this book contains a vocabulary test that can be used to measure whether the high-frequency words have been learned, and the progress of the learner in the learning of low-frequency vocabulary. *The Vocabulary Levels Test* exists in two different versions. There are also productive versions of the test (Laufer and Nation, 1995; Laufer and Nation, 1999) (see appendix 4). See Schmitt, Schmitt and Clapham (in press) for some research on this test.

The test is designed to be quick to take, easy to mark and easy to interpret. It gives credit for partial knowledge of words. Its main purpose is to let teachers quickly find out whether learners need to be

working on high-frequency or low-frequency words, and roughly how much work needs to be done on these words.

There is much more to vocabulary testing than simply testing if a learner can choose an appropriate meaning for a given word form, and we will look closely at testing in a later chapter. However, for the purpose of helping a teacher decide what kind of vocabulary work learners need to do, the levels test is reliable, valid and very practical.