

Word Meaning in Academic English: Homography in the Academic Word List

KAREN WANG MING-TZU and PAUL NATION

LALS, Victoria University of Wellington

The Academic Word List (Coxhead 2000) consists of 570 word families that are frequent and wide ranging in academic texts. It was created by counting the frequency, range, and evenness of spread of word forms in a specially constructed academic corpus. This study examines the words in the Academic Word List (AWL) to see if the existence of unrelated meanings for the same word form (homographs) has resulted in the inclusion of words in the list which would not be there if their clearly different meanings were distinguished. The study shows that only a small proportion of the word families contain homographs, and in almost all cases, one of the members of a pair or group of homographs is much more frequent and widely used than the others. Only three word families (*intelligence*, *offset*, and *panel*) drop out of the list because none of their homographs separately meet the criteria for inclusion in the list. A list of homographs in the AWL is provided, with frequencies for those where each of the members of a homograph pair are reasonably frequent.

For over seventy years (Dresher 1934), researchers have noted the occurrence of the same words across a range of quite different academic disciplines. These words, variously called ‘sub-technical vocabulary’ (Dresher 1934; Yang 1986; Flowerdew 1993), ‘semi-technical vocabulary’ (Farrell 1990), ‘academic vocabulary’ (Martin 1976; Coxhead 2000), ‘frame words’ (Higgins 1966), and ‘specialised non-technical lexis’ (Cohen *et al.* 1988), are distinct from the general high frequency words of the language. Typically, they are described as part of a set of three or four levels of vocabulary in academic text, starting with the 2000 word or so high frequency, general service vocabulary (*the*, *across*, *capital*, *discover*, *measure*, *small*, *wealth*), the academic vocabulary (*administration*, *assist*, *complex*, *economy*, *individual*, *minor*), the technical vocabulary closely associated with a particular subject area (*cursor*, *reboot*, *modem*, *rom*, *pixel*, *mouse* in computing), and low frequency vocabulary (*odoriferous*, *obliteration*, *obstreperous*, *panjandrum*, *panache*, *panacea*). Figure 1 consists of a text with each of these four types of vocabulary indicated.

The classic list of high frequency general service words is Michael West’s (1953) 2000-word *A General Service List of English Words*. This list is old and in need of replacement, but still works well (Nation and Hwang 1995), providing around 80 per cent coverage of the running words in academic texts.

The best researched and most recent list of academic words is Coxhead’s (2000) 570-word Academic Word List (AWL). This list covers at least 8.5 per

An *explication* is not simply a **definition**, for the only restriction on **definitions** is that the *abbreviating device* be *eliminable* in all places in favour of what it *abbreviates*. An *explicatum*, however, should behave, in certain *core areas*, in the same way that the *explicandum* behaves. In other words, it should accord with clear cut **intuitive judgements** that are made in certain **obvious** cases. For example, consider the *truth preserving inference*: it is raining, therefore it is raining. (from Oddie, G. 1986. *Likeness to Truth*. Dordrecht, Holland: D. Reidel Publishing Company, p. 3)

Unmarked = 1st 2000, **bold** = academic vocabulary, *italics* = low frequency, *underlined* = technical. Underlined items can be from the first 2000, AWL or low frequency words.

Figure 1: A text with different types of vocabulary marked

cent of the running words in a wide variety of academic texts. This academic vocabulary is characterized by (1) being reasonably frequent in most academic texts from a wide range of academic disciplines, (2) being relatively infrequent in other types of texts such as novels or colloquial spoken texts, (3) coming largely from French, Latin or Greek, and (4) being not obviously connected with any one particular subject area.

Coxhead's AWL was created from a corpus of 3,600,000 running words, called the Academic Corpus, from 28 subject areas equally distributed over the four divisions of Commerce, Law, Science, and Arts. Table 1 shows the composition of the corpus.

According to Coxhead, the AWL is made up of word families. Each family consists of a headword and its closely related inflected and derived forms according to level 6 of Bauer and Nation (1993), which includes all the inflections and a substantial number of derivational affixes. However, each derived form must contain a stem that can stand as a word in its own right (a free form). The meaning of the stem should not differ markedly between one member of the family and another. Table 2 contains some examples of families. However, when the family lists are used in computer programs to count vocabulary, because the counting is done on the basis of the word form, occurrences of homographs are counted as part of the same form.

Coxhead (2000) decided that to be included in the AWL, a word family had to (1) occur at least 100 times in the corpus (most occurred many more times than this), (2) occur in all four divisions of Commerce, Law, Science and Arts with a frequency of at least 10 in each division, and (3) occur in at least 15 of the 28 subject areas (most occurred in many more than this). Counting across the divisions and subject areas was done not on the assumption that subjects in each division shared common discourse practices, but simply because the list had to be useful for learners no matter what division or subject area they studied in. The major aim in making the list was to provide an explicitly described, feasible, vocabulary learning goal for learners with academic

Table 1: Subject areas in the faculty divisions of the Academic Corpus

Arts	Commerce	Law	Science
Education	Accounting	Constitutional law	Biology
History	Economics	Criminal law	Chemistry
Linguistics	Finance	Family law and medico-legal	Computer science
Philosophy	Industrial relations	International law	Geography
Politics	Management	Pure commercial law	Geology
Psychology	Marketing	Quasi-commercial law	Mathematics
Sociology	Public policy	Rights and remedies	Physics

Table 2: Sample families

A	
	AN
ABLE	
	ABILITY
	ABLER
	ABLEST
	ABLY
	ABILITIES
	UNABLE
	INABILITY
ACCOUNT	
	ACCOUNTED
	ACCOUNTING
	ACCOUNTS
	ACCOUNTANT
	ACCOUNTANTS
	ACCOUNTANCY

purposes. The list could be used for direct learning and teaching and in the choice and design of teaching materials.

The list is divided into 10 sublists, with sublist 1 containing the most frequent, wide-ranging words, and sublist 2 the next, and so on. Sublists 1 to 9 each contain 60 word families. Sublist 10 contains 30. Note that the unit of

counting in constructing the AWL is the word family. There is increasing research evidence showing that word families are psychologically real (Bertram *et al.* 2000; Nagy *et al.* 1989). Each family consists of a collection of related word forms, whether or not each form is found in the corpus used. For example, the *approach* family includes *approach*, *approachable*, *approached*, *approaches*, *approaching*, *unapproachable*, even though some may not have occurred in the academic corpus that was used to make the list, and even though the frequency of only one of the family members may have been sufficient to get the whole family into the list.

The AWL was largely created using a computer program called RANGE (available free on Paul Nation's web page at <http://www.vuw.ac.nz/lals/>) which can count the occurrence of words across many texts and can make use of pre-determined word families. A limitation of this computer-based approach is that the program counts word forms with no consideration of the meanings of these forms. Counting forms can misrepresent the composition and learning burden of a word family. For example the word family with the headword *row* is included in West's *General Service List*. However, the entry includes *row* (noun) a row of houses, *row* (verb) row a boat, and *row* (noun) = quarrel (see Figure 2). What is listed as one word family clearly should be three different families.

Moreover, if the quite different meanings of the form were considered, and if frequency figures alone were the criterion for inclusion in the list, none of these three meanings of *row* would have been included. Together, however, they are frequent enough for *row* to be included.

Does the AWL contain similar problem cases? For example, does the word *issue* get into the list because it has several different meanings (1) an important topic or problem ('We discussed the issue'), (2) the action of supplying or distributing ('We issued a new product'), (3) a publication ('He had an article in the latest issue of the journal'), (4) children or off-spring ('They died without issue'), and (5) the action of flowing out ('A stream of abuse issued from his mouth'), each of which might be common in different subject areas like politics, marketing, history, family law, and medicine? It is possible that the apparent wide-ranging occurrence of

ROW	174e		
row, n.	A.	A row of people, houses, seats, plants Standing in a row, in a row Second seat in the third row	62%
row, v.	B.	Row a boat Row a race [n. = <i>quarrel</i> , 10%]	25%

Figure 2: The entry for *row* from the *GSL* (West 1953: 420)

academic word families across various subject areas is an illusion created by counting word forms and ignoring their different meanings in the different subject areas.

So, simply counting forms creates a problem that may have affected the validity of the AWL. It is possible that the problem has been increased by grouping these forms together into word families. For example, the word family *consist* contains eleven family members *consist, consisted, consisting, consists, consistent* (which is the most frequent member in academic text), *consistently, consistency, consistencies, inconsistent, inconsistency, and inconsistencies*. Notice how *consistent* differs in meaning from *consist*, and how some uses of *consistency* (the consistency of glue—this meaning did not occur in the Academic Corpus) can differ from *consistent* (consistent behaviour). Similarly, the word family *project* has two meanings, 'a piece of work', as in 'the Genome project', and 'to go forward', as in 'a projection of future losses' and 'projected losses'. The first meaning occurred exclusively in the forms *project* and *projects*. The second meaning almost never occurred in the forms *project* and *projects* in the Academic corpus and most occurrences were in the forms *projected, projecting, projections, projector, and projectile*. Thus, it could be a misrepresentation to say that the Academic Word List contains 570 word families that would all be very useful for learners with academic purposes no matter what subject they are studying, and the principal goal of this study is to investigate whether counting word forms without concern for homography has resulted in a list of academic words that differ substantially in their meaning from one discipline to another.

So far we have talked about different meanings of a word form as if these different meanings were unrelated to each other. But this is not always so. The word *conceive*, for example, can be used to talk about creating an idea ('We could not conceive of a way of dealing with this problem') or creating a baby ('The child was conceived out of wedlock'). These two uses have many similarities, and rather than consider them as different words, they could be considered as closely related different uses of the same word. Two or more clearly related meanings of the same form are called *polysemes* and the phenomenon is called *polysemy*. Two or more completely unrelated meanings of the same written form (*row* of houses, and *row* a boat) are called *homographs* and the phenomenon *homography*. Polysemy and homography are points on a scale and there can be considerable disagreement about whether two items are polysemes or homographs. Some dictionary makers tend to distinguish meanings and thus move towards the homography end of the scale. Language teachers tend to show how meanings are related and how knowing one can help understand the other and thus they tend to move towards the polysemy end of the scale when dealing with different uses of a word. This can be seen in the numerous related senses listed for words in some dictionaries. Thus a major problem in dealing with the meanings of words is to decide when the various senses are instances of polysemy or when they are homographs. Because we are dealing with items

in a written corpus, the term *homograph* is used instead of the more inclusive term *homonym* which can include spoken or written similarity. Some items in the study, such as *attribute* (noun) and *attribute* (verb), *converse* (adjective) and *converse* (verb), and *abstract* (adjective and noun) and *abstract* (verb), are homographs and not homophones.

With these issues in mind, the present study set out to answer the following research questions.

- 1 What proportion of the word families in the AWL contain homographs that occur frequently in the Academic Corpus?
- 2 If these homographs are counted as separate word families, what word families in the AWL would not meet Coxhead's criteria for inclusion in the AWL?
- 3 If a word family contains homographs, to what degree are the different homographs limited to one particular discipline or larger subject area?

To answer these questions it is necessary to be able to distinguish polysemic and homographic uses of each word form and each word family.

METHOD

The main task faced in this research is to determine (1) which words in the AWL contain homographs, and (2) which of these homographs are frequent enough in the academic corpus to affect the inclusion of the word family in the AWL. The major steps in the investigation are

- 1 to examine the various meanings of each member of each AWL word family to find which families contained homographs, that is which ones were really two or more different families. A dictionary and a rating scale for relatedness of word meanings were used for this purpose.
- 2 to see how frequent and widely distributed each homograph is in order to find out if each homograph still meets the criteria for inclusion in the AWL. This was done by using a concordance program on the Academic Corpus and tabulating the results. The AWL criteria were applied to the results.

It was decided to proceed by taking representative samples from the AWL with the eventual aim of covering the whole AWL. The first sample consisted of one-fifth of the list, that is 114 word families, so that the results in the first sample could be used to establish a hypothesis and help to shape the design of the procedure to be used in the collection and analysis of further samples. The word families in the first sample were chosen by examining every fifth word in the sublists. A dictionary was used to locate words with homographs. After a systematic and detailed examination of several dictionaries, the *New Oxford Dictionary of English* (Pearsall 1998) was chosen for this stage of the analysis. This dictionary is recent, has a large number of entries, includes academic uses of words, and does not contain a lot of archaic usage.

A semantic relatedness scale was developed to analyse the dictionary definitions to find the word families of the AWL that contained homographs. These were then used as a basis for the gathering and analysis of concordance data. Eventually, all 570 items in the list were investigated. An inter-rater reliability check was carried out to test the reliability of applying the scale.

The design of a semantic relatedness scale

A semantic relatedness scale was designed to group the definitions given in the dictionary in order to reliably determine whether a word family contained homographs (that is the definitions were unrelated to each other), or if all the definitions for the word forms were sufficiently closely related to be considered as different senses or polysemes. The semantic relatedness scale (Table 3) developed in this study was based on the one used by Nagy and Anderson (1984). The scale has six categories, from 'Level 0' to 'Level 5'. Two meanings are classified as related at Level 0 if their meanings are the same. For example, the level of relatedness of meanings of the word *intermediate* in '*intermediate technology*' and '*intermediated between the individual and the state*' is classified as Level 0, as they both have the base or shared meaning of 'coming between two things'. In this example, the differences in parts of speech were ignored as the base form '*intermediate*' has the same base meaning in both types. The 'base' meaning is the common meaning shared by two or more senses of the word forms. Using this scale, the relatedness between every two base meanings is first placed at one of these six levels and then categorized as polysemy or homography. The cut-off point between polysemes and

Table 3: *The semantic relatedness scale*

Semantic relatedness level	Description of the degree of relatedness
0	The meaning is the same as the base meaning.
1	The meaning is only slightly different from the base meaning.
2	The meaning is related to the base meaning with some changes.
3	The meaning is substantially different from but is still related to the base meaning.
4	The meaning is very distantly related and almost totally different from the base meaning.
5	There is no relationship at all between this meaning and the base meaning.

homographs is between Level 3 and Level 4. This means that all of the base meanings in the same group should be related at Level 3 or lower.

The rating scale in Table 3 is an ordinal scale representing the ranking of relatedness. Each of the levels is not necessarily of the same 'size' or 'length'. For example, the following definitions for *promotion* are from the *New Oxford Dictionary of English*:

- 1 activity that supports or provides active encouragement for the furtherance of a cause, a venture, or aim, e.g. GPs have a vital role to play in health promotion;
- 2 a sporting event, especially a series of boxing matches, staged for profit, e.g. a boxing promotion;
- 3 the action of raising someone to a higher position or rank or the fact of being so raised, e.g. majors designated for promotion to lieutenant colonel.

The base meaning of *promotion* is something like 'putting something or someone forward or up so that it gains more importance or notice in the eyes of others'. The third definition, 'the action of raising someone . . .' comes closest to this. The second definition, 'a sporting event . . .' is closely related to definition 1 'furtherance of a cause'. Definitions 1 and 2 are both related to definition 3 with some differences that have been assessed to be Level 2 differences (see Table 3). However, it is difficult to determine if these differences are of the same 'size', but this is not important for this study.

The semantic relatedness scale was designed using the following two principles:

1 The current study adopts Ruhl's (1989) framework of 'Monosemic Bias' and assumes that there is an abstract general meaning for all senses of a word. For example, in all the following senses of the word *neutral*, there is clearly a shared meaning: (1) referring to a country that does not officially support either side in a war, (2) referring to a position not supporting either side in an argument, (3) not showing emotions or preference, (4) referring to a voice not showing emotions, (5) not causing change because it is equal on both sides, (6) a position between the gears of a car, (7) referring to electrical wires—neither earth nor positive, (8) referring to colours, (9) referring to atomic particles—neither positive nor negative, (10) neither acid nor alkali.

If such a general meaning cannot be found, then two or more new general meanings need to be created to represent the homographs. If the meanings are almost totally different and only very distantly related (Level 4 relatedness), the connection is probably only discernible for highly competent speakers with a good knowledge of a particular word. It would not be reasonable to expect less proficient learners to be able to predict or guess a meaning from a very distantly related meaning. Such relatedness is categorized as one between homographs.

2 The procedure is designed and used in such a way that the results from the meaning analysis are dictated by the dictionary entries as little as possible.

For example, different senses with seemingly related definitions are not necessarily included under one general meaning when the scale has been applied. Thus, there is a vague connection between *issue* meaning 'to flow out' and *issue* meaning 'child' as used in legal documents, but they were put into different word families even though the dictionary lists them under one entry. Thus the dictionary senses were used as a guide but they were not adopted uncritically and were combined or separated when necessary.

Instead of using a scale, it might have been possible to look for systematic relationships between senses. Polysemy makes use of such systematic relationships within a conceptual area. These systematic relationships can be used to produce and comprehend new senses. The mechanisms used to do this include metaphor, metonymy, and conventional knowledge. The main goal of the scale, however, is to distinguish homographs. It does this by helping the rater to use intuitive judgement to rate the level of observable semantic relationship between two uses of the same word form. Insights from cognitive semantics (Lakoff 1987) could be used to explain the relationship between a base meaning and its related senses, but a separate judgement would still need to be made regarding the level of closeness of this relationship. The scale indicates that homonymy is not necessarily an arbitrary or coincidental distinction although some researchers in cognitive semantics claim it is. What was once a transparent systematic relationship may, over time, become opaque or tenuous, resulting in a shift from polysemy to homonymy.

Using the scale

The guidelines for the application of the scale are as follows:

- 1 When a sense is related to several basic general meanings, it should be grouped with the sense to which it is most closely related.
- 2 The final combination of general base meanings should be the one that makes as few meaning groups as possible.
- 3 Each sense in the same base meaning group has to be related to the base meaning but does not necessarily have to be related to all of the other senses in the same group. This can be explained by Wittgenstein's (1953) theory of 'Family Resemblances' which states that 'members of a category may be related to one another without all members having any properties in common that define the category' but this is using the term *category* in a different way from which the term *family* is used in this study.
- 4 The grouping of meanings is a process involving a series of iterations. The result with the least number of base meaning groups might not be obtained until after several iterations.
- 5 The key phrases chosen for the base meanings in the final results should best describe the 'family resemblance' of that meaning group.
- 6 In borderline cases, where a sense is related to two base meanings at the

same level, it should be included in the stronger group, that is, the group whose general meaning comprises more senses.

- 7 Archaic uses and proper nouns were not included in the analysis, e.g. *factor* meaning 'business agent', *Bond* as in *James Bond*.

Thus, the method of analysis involved going through all the senses of a word family in the dictionary, and grouping them using the above guidelines. This was done starting with each of the meanings and considering all the others each time. The aim of the analysis was to make the smallest number of different word families (homographs), with the base meanings contained in any one word family as closely related to each other as possible. Let us look at an example using the word family *volume*. The base meanings found after the initial analysis of the dictionary entry are:

- 1 a book;
- 2 a consecutive sequence of issues of a periodical;
- 3 the amount of space;
- 4 an amount/quantity;
- 5 power of sound.

These meanings are listed across the top of the table. The numbers followed by a closing bracket, for example 2), placed within the table refer to these meanings. The relatedness between each two of these base meanings is summarized in Table 4.

Column 2 of the table shows that meaning 1) *a book* is related to meaning 2) at Level 2, to meanings 3) and 4) at Level 4, and to meaning 5) at Level 5.

Initial analysis showed that there are two main meaning groups: 1) and 2) are more closely related to each other than they are to 3), 4), and 5). As the table shows, 1) and 2) are related at Level 2 of the scale (see columns 2 and 3 of Table 4), which is well below the cut-off point of Level 4. For each of the main groups, a base meaning description needs to be chosen as the phrase to describe each group. Let us look at these two groups in turn.

If 3) *the amount of space* (column 4) was chosen as the meaning description for one of the homographs, it would include 4) as one of its polysemous meanings, because it is at Level 3, but would have 1), 2), and 5) as three of its homographs, because they are beyond Level 3. Choosing 5) *the power of sound* (column 6) as the meaning description for the homograph gives the same result. To use 4) *an amount and quantity* (column 5) as the homograph meaning, the description would include 3) and 5) as its polysemous meanings, that is, it includes more of other base meanings than the first two choices. So 4) is the best choice. Note that even though meanings 3) and 5) are both related to 4) at Level 3 of the scale (as shown in column 5), 5) is related to 3) at Level 4 of the scale (as shown in column 6). So, choosing 4) *an amount and quantity* as the second base meaning is the best choice because meanings 5) *the power of sound* and 3) *the amount of space* are related to it within the first three levels of the scale which is the criterion necessary for the meanings to be polysemes.

Table 4: Scale of semantic relatedness for the word family *volume*. (*volume/volumes/vol.*) (Sublist 3)

Semantic relatedness level	1) a book	2) a consecutive sequence of issues of a periodical	3) the amount of space	4) an amount/ quantity	5) the power of sound
0					
1					
2	2) a consecutive sequence of issues of a periodical	1) a book			
3			4) an amount/ quantity	3) the amount of space; 5) the power of sound	4) an amount/ quantity
4	3) the amount of space; 4) an amount/ quantity		1) a book; 5) the power of sound	1) a book	3) the amount of space
5	5) the power of sound	3) the amount of space; 4) an amount/ quantity; 5) the power of sound	2) a consecutive sequence of issues of a periodical	2) a consecutive sequence of issues of a periodical	1) a book; 2) a consecutive sequence of issues of a periodical

Meanings 1) and 2) belong to the same main group. Meaning 1) *a book* is a better description of the meaning in this homograph group than 2) *a consecutive sequence of issues of a periodical*, because it is less specific and gives fewer related items at Level 5. Meanings 3), 4), and 5) belong to the second group, with 4) *an amount and quantity* being the best description of this group. All of the meanings in each group are related to that group's base meaning below Level 4 of the scale.

Note that although 1) and 2) are related to each other at Level 2, they relate to the other three base meanings at different levels: 1) relates to all of 3), 4), and 5) at Level 4 whereas 2) relates to 3), 4), and 5) at Level 5.

The results of the classification show that the word family *volume* has two unrelated base meanings: the primary meaning is related to ‘a book’ and the secondary meaning is ‘an amount and quantity’.

Here is another example using the family *issue* (*issue, issuing, issued, issues*) (Sublist 1). The base meanings found after an initial analysis are:

- 1 an important topic;
- 2 the action of distributing; a publication;
- 3 children;
- 4 to flow out;
- 5 a result or outcome.

Table 5: Scale of semantic relatedness for the word family *issue* (*issues/ issued/issues/issuing*) (Sublist 1)

Semantic relatedness level	1) an important topic	2) the action of distributing	3) children	4) the action of flowing	5) a result or outcome
0					
1				5) a result or outcome	4) the action of flowing
2					
3		4) the action of flowing		2) the action of distributing;	
4	4) the action of flowing; 5) a result or outcome	3) children; 5) a result or outcome	2) the action of distributing; 4) the action of flowing; 5) a result or outcome	1) an important topic; 3) children	1) an important topic; 2) the action of distributing; 3) children
5	2) the action of distributing; 3) children	1) an important topic	1) an important topic		

It is clear in Table 5 that base meaning 1) is very distant from all the other four senses, so the first homograph of *issue* is the sense of ‘an important topic’. Base meaning 4), ‘the action of flowing’ and base meaning 5) ‘a result or outcome’ are related at Level 1 so can be combined as one meaning group. This group was named as ‘the action of flowing’ because this is the most general description of the three base meanings. Base meaning 2) is related to base meaning 4) at Level 3, so is placed into the ‘the action of flowing’ base meaning group. Base meaning 3), ‘children’, is related to base meaning 4) at Level 4, so that should be in a separate meaning group. So these meanings are

then divided into three meaning groups: 1) an important topic; 2) flowing; 3) children.

Note that although base meanings 4) and 5) (columns 5 and 6 in Table 5) are related at Level 1, they both relate to base meaning 3) at Level 4, but relate to base meaning 2) at different levels. The goal of this analysis was to distinguish homographs and to make sure that each homograph contained senses that were closely related with each other.

This part of the analysis was carried out in a decontextualized way using dictionary entries. It was the most efficient way, because if the analysis was done on a concordance, the same kinds of generalizations present in the dictionary would have to be made in order to end up with groupings of senses. In the later stages of this study, the meaning groupings were checked against concordances when the frequency, range, and distribution data were being gathered.

Checking the use of the scale

An inter-rater reliability test was carried out to test if the use of the rating scale gives a reasonable level of reliability in categorizing the word families as with or without homographs. Because the scale requires judgements based on intuition, it is important to show that the classification can be done similarly by different individuals.

The first researcher was the original rater. The other rater, who was only used in the inter-rater check, was a native speaker of English who is a qualified and experienced ESOL teacher. The word families used in the inter-rater reliability check were chosen from the first sample of 114 word families. In the preparation, these were first divided into two categories: word families (1) with homographs, and (2) without homographs. These were then further divided into two groups according to the level of difficulty in coding them. The level of difficulty was mainly based on the complexity of the network of the meanings in a word family, using the time it took to reach a decision as a guideline. Word families that took more than five minutes to rate were classified as 'difficult'. In all, there were four different categories: word families with homographs that were easy to rate, families without homographs that were easy to rate, families with homographs that were difficult to rate, and families without homographs that were difficult to rate. After training, 20 word families, five from each of the four different categories, were arranged randomly for the rater to analyse independently. Each word family had all the members in that family listed. Dictionary entries to be analysed were marked. This number of items well exceeds the minimum number of five items needed to establish accuracy of coding involving two alternatives at the 0.05 significance level (Rosenthal 1987: 64). The second rater coded the 20 word families prepared. This took 50 minutes. A raw accuracy score of 85 per cent is desirable for coding involving two alternatives, although 70 per

cent is acceptable (Rosenthal 1987: 67). The raw accuracy score of the reliability check was 75 per cent as calculated below:

Number of agreements in coding word families with homographs = 7
(out of 10)

Number of agreements in coding word families without homographs = 8
(out of 10)

Number of agreements = 7 + 8 = 15, Total number of word families coded = 20

Accuracy score = $(7 + 8) / 20 = 75$ per cent

Although it would be better to have a higher level of reliability, this is a satisfactory level, although it is clear that there is a fair degree of subjectivity involved in the decision making. The inter-rater check did not look at the various levels of relatedness.

The identification of homographs using the dictionary had the goals of (1) identifying the word families in the AWL that contained homographs and thus whose family members needed to be looked at further using concordance data, and (2) providing a set of meanings to use when classifying the concordance data of the family members to obtain frequency data for each meaning.

The concordance data

The WordSmith Tools program (Scott 1999) was used to search the 3,500,000 token Academic Corpus used in the original AWL study for all of the different word types of the AWL word families that had been found to contain homographs after analysing the dictionary definitions. Each item found appeared in a context set at 10 words to the left and 10 to the right although it was extended to as many as 50 words to the left and 50 to the right when it was more difficult to determine which of the meaning groups was used.

Word families from Sublist 1 have up to 3,118 entries in the total corpus (the word family *policy*) and those from Sublist 2 have up to 993 entries (the word family *normal*). These are much higher than those in the other sublists so it was decided that samples should be taken for word families from Sublist 1 and Sublist 2. This was done so that there would not be enormous numbers of concordance items to classify. For the concordance samples taken to be representative of the respective populations, the required minimum sample size for a population of 3,500 (which exceeds the 3,118 of the largest number of occurrences) is 346 and for a population of 1,000 is 278 (Krejcie and Morgan 1970: 608). As shown below, the minimum population sizes were always exceeded.

Every fourth entry was selected to compile a file of about 400 entries for word families from Sublist 1 and every third entry was selected to compile a file of about 350 entries for word families from Sublist 2. The smallest number of entries sampled for Sublist 1 families was 388. This was for the word family *policy*, which has a total frequency of 3,118. *Normal*, which has a family

frequency of 993, had a sample size of 337, which is the least number of entries taken for a Sublist 2 word family. All of the sample sizes were greater than the minimum required for the samples to be statistically representative of their respective family frequencies in the corpus.

For sublists 3 to 10, the complete set of concordances from the Academic Corpus was examined for each word family. After the first sample of 114 families was completed, a second sample was analysed using concordance data.

RESULTS AND DISCUSSION

Of the 570 word families in the AWL, 60, approximately 10 per cent, contained words which were actually homographs and that therefore should not be listed under the same family. These words are listed in the Appendix and were the focus of the concordance study. Only 21 of these homographs occurred frequently enough in the Academic Corpus to be worth further consideration. These are listed in Table 6. They are listed separately because the homographs were frequent enough to affect the status of the family in the AWL, and to be of interest when these words are dealt with in class. The remaining 39 word families contained homographs that did not occur at all in the Academic Corpus, or occurred with a frequency of less than 5 per cent of the frequency of the total word family. For example, *convert* meaning 'to convert a try in rugby', and *function* meaning 'a formal event or ceremony' did not occur at all. *Proceeds* meaning 'money obtained from an event or activity' and *appropriate* meaning 'to take for one's use, usually without the owner's permission' make up less than 5 per cent of the occurrences of each word family.

For these 21 word families in the AWL that are to be considered in further detail, additional families had to be created because in many cases some family members and sometimes all family members had homographs. These additional families were created so that their frequency, range, and distribution could be examined in the academic corpus to see if they met the criteria for inclusion in the AWL, namely a total frequency of more than 100, and a frequency of at least 10 in each of the four divisions of the corpus. Six of these 21 additional word families did meet the criteria: *consist* (made up of); *issue* (the action of flowing); *volume* (book); *attribute* (feature); *objective* (not subjective); and *abstract* (extract, remove as in the abstract of an article). In addition, as Table 6 shows, only three of these 21 word families, *intelligence*, *offset*, and *panel*, drop out of the AWL completely, because neither the additional families created because of homography nor the remaining reduced original family met the AWL criteria. These word families were already quite low in the AWL sublists, namely in sublists 6, 8, and 10. So, of the total 60 families with homographs, 21 had homographs that were frequent enough in the academic corpus to get further attention, and of these 21, 6 had members which met the criteria for being additional AWL families, and 3 of the 21 had

to be excluded from the AWL because neither homograph in each family met the criteria for inclusion. For the remaining 12 of the 21, no change was required in the AWL because the homographs were so infrequent. So for the 570 original word families, there were only 9 changes (6 additions and 3 deletions).

Table 6 provides the details of the occurrences of the meanings of the 21 word families. The word family *consist*, for example, occurred a total of 1,549 times in the academic corpus, 1,089 times (70 per cent of occurrences) meaning 'staying the same', and 460 times (30 per cent of total occurrences) meaning 'made up of'. Both uses should be separately listed in the AWL because, as Table 6 shows, they both occur frequently across the four sub-corpora of Arts, Commerce, Law, and Science. Note that the meanings that are excluded (shown in brackets) either had a frequency of less than 100, as in *issue* M3, or did not occur or had fewer than ten occurrences in one or more of the sub-corpora, as in *issue* M3, *credit* M2. Note also how none of the meanings of *intelligent*, *offset*, and *panel* meet the frequency or range criteria.

The results of this study counter the criticism that the AWL is the result of counting the same forms with unrelated meanings using word families as the unit of counting. Homography played a very minor role in determining what was included in the AWL.

To a certain degree, the use of word families inflates the problems caused by homography. That is, if word types or even lemmas were listed instead of families, the amount of homography would be less because different meanings tend to be represented by different types. For example, *Orient* meaning 'the East' does not occur in the forms *oriented*, *orienting*, and *orients*. *Orient* meaning to face the right direction does not occur in the form *oriental* and does not need a capital as *Orient* (the East) does. So, within a word family based on form there may be formal differences between homographs. There are numerous examples of this in the words investigated in this study, for example *proceeds* meaning 'money obtained from an event', *conversation* as a part of the *converse* word family, and *tense* meaning 'a verb form', which never has the form *tension* as a member of the *tense* word family. In addition, there are plenty of grammatical and collocation clues to distinguish homographs. *Consist* with meaning 1 (staying the same) has the forms *consistency*, *consistencies*, *inconsistencies*, *consistent*, *inconsistent*, and *consistently*. *Consist* with meaning 2 (made up of) has the forms (all verb based) *consist*, *consists*, *consisted*, and *consisting*. *Attribute* with meaning 1 (to ascribe, accredit) has the forms (all verb based) *attributable*, *attributed*, *attributes*, *attributing*, *attribution(s)*, *attributational*. With meaning 2 (a feature), *attribute* has the noun forms *attribute(s)*. There are many additional examples of this formal differentiation of meanings, a point frequently noted by Sinclair (1991). The choice between (1) types or lemmas, and (2) word families as the unit of counting is a critical one in corpus-based studies of vocabulary. For making up the AWL, word families was the best choice because learners dealing with the words in the AWL would already have some knowledge of English affixes, and, as this

Table 6: Distribution of primary and other meanings of the 21 word families with significantly frequent homographs

			Frequency in each of the sub-corpora					
			Total	%	Arts	Comm	Law	Science
consist			1549		342	421	234	552
Sublist 1	M1	{consistent} staying the same	1089	70	196	367	164	362
	M2	made up of	460	30	146	54	70	190
issue			2701		586	1071	873	171
Sublist 1	M1	an important topic	1982	73	*445	*876	*564	*97
	M2	the action of flowing or producing children)	699	26	*141	*195	*289	*74
	(M3)		20	1	*0	*0	*20	*0
credit			1194		98	346	719	31
Sublist 2	M1	a payment received etc to acknowledge or approve)	1097	92	11	344	711	31
	(M2)		97	8	87	2	8	0
normal			993		193	195	187	418
Sublist 2	M1	relating to a standard perpendicular)	899	91	*193	*195	*187	*324
	(M2)		94	9	*0	*0	*0	*94
correspond			679		114	115	84	366
Sublist 3	M1	to match or agree	602	89	82	103	54	363
	(M2)	to write)	77	11	32	12	30	3
volume			711		105	157	100	349
Sublist 3	M1	a quantity	361	51	22	56	21	262
	M2	a book	350	49	83	101	79	87
attribute			528		109	141	177	101
Sublist 4	M1	to ascribe; to accredit	432	82	99	128	161	44
	M2	a feature	96	18	10	13	16	57
project			579		125	178	65	211
Sublist 4	M1	a piece of work	435	75	77	158	60	140
	(M2)	to go forward, stick out, predict)	144	25	48	20	5	71
decline			438		179	136	52	71
Sublist 5	M1	to decrease	381	87	177	123	10	71
	(M2)	to formally refuse)	57	13	2	13	42	0
generation			477		127	83	36	231
Sublist 5	M1	a stage of development production)	350	73	121	22	35	172
	(M2)		127	27	6	61	1	59
objective			978		265	454	144	115
Sublist 5	M1	a goal	752	77	212	366	79	95

Table 6: cont.

				Frequency in each of the sub-corpora					
				Total	%	Arts	Comm	Law	Science
abstract				375		113	91	32	139
Sublist 6	M1	to extract, remove, summarize		186	50	46	57	13	70
	M2	not concrete		189	50	67	34	19	69
	(M3)	{abstracted} preoccupied)		0	0	0	0	0	0
brief				337		116	77	65	79
Sublist 6	M1	short		301	89	104	59	59	79
	(M2)	instructions)		36	11	12	18	6	0
intelligent				331		52	196	11	72
Sublist 6	(M1)	able to acquire knowledge)		141	43	49	14	6	72
	(M2)	{intelligence} the collection of information)		190	57	3	182	5	0
appreciate				202		*50	64	43	45
Sublist 8	M1	to value, understand		177	88	48	52	38	39
	(M2)	to increase in monetary value)		25	12	2	12	5	6
offset				221		10	77	40	94
Sublist 8	(M1)	counteract)		120	55	0	77	40	3
	(M2)	place out of line)		96	43	5	0	0	91
	(M3)	beginning)		5	2	5	0	0	0
tense				271		120	47	25	79
Sublist 8	M1	rigid		217	80	67	47	24	79
	(M2)	a verb form)		54	20	53	0	1	0
induce				262		34	72	31	125
Sublist 8	M1	to cause		235	90	27	71	30	107
	(M2)	a way of reasoning)		27	10	7	1	1	18
accommodate				177		58	57	40	22
Sublist 9	M1	to adapt		127	72	47	36	25	19
	(M2)	to provide lodging)		50	28	11	21	15	3
converse				308		182	40	48	38
Sublist 9	M1	the reverse		128	42	24	33	35	36
	(M2)	to have a conversation)		180	58	158	7	13	2
panel				121		18	60	11	32
Sublist 10	(M1)	a group of people)		61	50	9	33	7	10
	(M2)	a sheet of material)		60	50	9	27	4	22

M1 = primary meaning M2 = secondary meaning M3 = tertiary meaning

* = estimated figures, () indicates a meaning not meeting the criteria for inclusion in the AWL. Where the head word of the new family is different in form from the original stem, this is shown in wavy brackets {}.

study shows, the homography problems that occur because of word families affect only a small number and a small proportion of the word families.

This study set out to answer three questions regarding (1) the number of families containing homographs in the AWL, (2) the changes that would need to be made to the AWL to take account of these homographs, and (3) the degree to which particular meanings of a word are limited to a particular subject area.

It was found that 60 families contained potential homographs. Only 21 of these actually contained homographs which occurred in the Academic Corpus, and only six of these would require additional entries in the AWL. Three word families would need to be removed from the AWL. These are reassuringly small changes, indicating that homography is not a major factor affecting the words in the AWL. Where homographs do occur, one of the homographs is typically much more frequent than the other, accounting for at least 95 per cent of the combined frequencies of the homographs. Just as some words occur much more frequently than other words, some meanings of a word occur much more frequently than others (Table 6 and Appendix 1 provide information about which meanings are the most frequent). This finding probably applies to most form-based studies. While it should not be taken as licence to ignore meaning when doing computer-based counting, it shows that the results of such counting are unlikely to be strongly affected by homography. The study does suggest that the unit of counting (word family, lemma, word type) needs careful consideration.

There are some word meanings that are limited to particular subject areas (see Table 6). Twenty of the 21 occurrences of *issue* meaning 'children' occur in Law; all 9 of the occurrences of *normal* meaning 'perpendicular' occur in Science; 182 of the 190 occurrences of *intelligence* meaning 'collection of information' occur in Commerce. These items, however, do not meet the frequency, range, and distribution criteria for inclusion in the AWL. Specialist meanings did not have a significant effect on the inclusion of word families in the AWL.

There are a few cases where a family included in the AWL has few occurrences in a particular subject area. *Decline* (to decrease) has only 10 occurrences in Arts compared with 71 up to 177 in the other areas. *Credit* (a payment received) has only 11 occurrences in Arts compared with 31 up to 711 in other areas. But these are not typical. In general, it is possible to say that the words in the AWL do not significantly change their meaning when they occur in the different sub-divisions of the Academic Corpus. This has implications for teaching.

Implications for teaching and learning

One of the assumptions behind this study, the monosemic bias, is that for receptive purposes polysemes should be treated as having a common underlying meaning. The teaching and study of words in the AWL should

reflect this. That is, learners should be encouraged to look for the central concept behind a variety of uses. This is important not only for academic vocabulary, but also for technical vocabulary where known words, like *price*, *cost*, and *demand*, take on technical meanings in a field like Economics. There are several related ways of drawing attention to the shared features of polysemes. One way is to do what was done in the early stages of this study, to look at dictionary entries for a polysemic word and relate the sub-entries to each other (Table 7). This can be done by trying to create a unifying definition, by selecting a core use, or diagrammatically by connecting the meanings to each other in some form of visual representation. Visser (1989) devised exercises to do this (see Table 7). The exercises can be done individually or with learners working in pairs.

Another way is to do what was done in the later stages of this study and to work with learners to classify items in concordances looking for formal clues to the different meanings, such as the particular word type and collocations, and looking for systematic relationships. These systematic and predictable relationships can involve figurative devices like personification, and states becoming actions and vice versa. Cognitive semantics and its treatment of polysemy can provide very helpful guidelines here. At the very least,

Table 7

Interpret: verb	Interpret: verb	What is the core meaning of this word?
<p>If you interpret something in a particular way, you decide that this is its meaning or significance. <i>Even so, the move was interpreted as a defeat for Mr Gorbachev . . . The judge says that he has to interpret the law as it's been passed . . . Both of them agree on what is in the poem, but not on how it should be interpreted.</i></p> <p>How would you interpret the meaning of this sign? </p>	<p>If you interpret what someone is saying, you translate it immediately into another language.</p> <p><i>The woman spoke little English, so her husband came with her to interpret . . . Three interpreters looked over the text for about three or four hours and found that they could not interpret half of it.</i></p> <p>Which is usually more difficult—interpreting from English to your first language, or interpreting from your first language to English?</p>	

understanding of the devices behind polysemy can help the learning and retention of such senses (Kövecses and Szabo 1996). More ambitiously, it can help in the interpretation and use of polysemes. The focus of discussion would be directed towards these devices rather than the particular words, that is, towards the systems rather than the outcomes.

The skill of interpreting polysemes is likely to be closely related to skill in guessing meanings from context, and a third way of focusing on polysemes would be to draw attention to polysemic words when they occur in intensive reading, getting the learners to recall related uses, and discussing their relationships.

This study can be seen as one kind of validation of the AWL and as an attempt to bring about some small improvements in the list and its use. It also sheds some light on the nature of academic vocabulary and the nature of word families, particularly as a unit of word counting. In addition, it provides information on the empirical basis for the concepts of polysemy and homonymy. Finally, the study is an example of how looking at what to learn and teach can also provide guidelines on how to learn and teach.

(Final version received January 2004)

APPENDIX: THE SIXTY AWL WORD FAMILIES CONTAINING HOMOGRAPHS

Infrequent meanings (less than 5 per cent of total occurrences) are in brackets. Families in italics have frequent homographs and can be found in Table 6. Where the head word of the new family is different in form from the original stem, this is shown in wavy brackets.

Sublist 1

assume

to accept as true without proof
(to take on [a role])

consist

made up of
{consistent} *staying the same*

contract

a written agreement
(to decrease in size)
(to catch a disease)

function

a purpose
(a formal event)

issue

an important topic

the action of flowing or producing
(children)

major

important
(a rank in the army)

period

a period of time
(a full stop)

policy

a course or principle of action
(an insurance contract)

proceed

to begin a course of action
{proceeds} (money obtained from an event)

Sublist 2

affect
 to have an effect on
 (to pretend to have or feel something)
 {affective} (emotion and desire)

appropriate
 suitable or proper
 (to take for one's use)

chapter
 a main division of a book
 (a branch of an association)

commission
 an official group of people
 (particular work for payment)
 (to perform)

credit
a payment received
to acknowledge or approve

institute
 an organization
 (to cause change)

normal
conforming to a standard
perpendicular
([a school] used for teacher training)

purchase
 to buy
 (firm contact)

Sublist 3

convene
 (to call for a meeting)
 {convention} way in which something
 is done

correspond
to match or agree
to communicate by mail

deduction
 an amount taken away
 (an inference)

demonstrate
 to show
 (to protest)

justify
 to prove
 (to line up print)

layer
 a sheet of material
 (a person or thing that lays something)

minor
 less important
 (children)

remove
 to take away
 (the degree of remoteness)

volume
a book
an amount, quantity

Sublist 4

attribute
to ascribe, accredit
a feature

project
a piece of work
to go forward, stick out, predict

promote
 to raise a position/publicity
 (to institute as a prosecution)

Sublist 5

compound
 mixture
 (to forbear from prosecution)
 (a fenced area)

decline
to decrease
to formally refuse

draft
 a plan
 (to enlist)

generation
stage of development
production

objective
a goal
not subjective

orient
 to determine relative position
 {the Orient} (East Asia)

prime
 of first importance
 (to make something ready for use)

stable
 firmly fixed
 (a place for horses)

Sublist 6

abstract
to extract, remove, summarise
not concrete
{abstracted} (preoccupied)
 attach
 to fasten, join
 (to seize by legal authority)
brief
short
instructions
intelligent
able to acquire knowledge
{intelligence} the collection of information

Sublist 7

convert
 to change
 (to score [in rugby])
 dispose
 to get rid of, place
 {disposed} (inclined)
 file
 folder holding documents
 (a line)
 (a tool for shaping)
 grade
 a level, slope
 (to cross livestock)
 media
 means of communication
 ((Anatomy) an intermediate layer)
 quote
 words taken from a text
 (an estimated cost)
 sole
 one and only
 (the bottom of the feet)
 (a fish)

Sublist 8

appreciate
to value, understand
to increase in value
 exhibit
 to show, display
 {an exhibition} (an award or scholarship)
 exploit
 to make full use
 (a bold feat)
induce
to cause
a way of reasoning
offset
to counteract
to place out of line
(to transfer an impression (printing))
 radical
 progressive
 (relating to the root)
tense
stretched tight or rigid
a verb form

Sublist 9

accommodate
to adapt
to provide lodging
 bulk
 a large amount
 (roughage in food)
converse
the reverse of something
(to have a conversation)
 norm
 a standard, a usual situation
 ((Maths) product of two conjugates)

Sublist 10

panel
a group of people in a discussion
a sheet of material

REFERENCES

- Bauer, L. and I. S. P. Nation. 1993. 'Word families.' *International Journal of Lexicography* 6: 253-79.
- Bertram, R., M. Laine, and M. Virkkala. 2000. 'The role of derivational morphology in vocabulary acquisition: Get by with a little help from my morpheme friends.' *Scandinavian Journal of Psychology* 41: 287-96.
- Cohen, A., H. Glasman, P.R. Rosenbaum-Cohen, J. Ferrara, and J. Fine. 1988. 'Reading English for specialised purposes: discourse analysis and the use of student informants' in P. Carrell, J. Devine and D. E. Eskey (eds): *Interactive Approaches to Second Language Reading*. Cambridge: Cambridge University Press. pp. 152-67.
- Coxhead, A. 2000. 'A new academic word list.' *TESOL Quarterly* 34: 213-38.
- Dresher, R. 1934. 'Training in mathematics vocabulary.' *Educational Research Bulletin* 13: 201-4.
- Farrell, P. 1990. *Vocabulary in ESP: a lexical analysis of the English of electronics and a study of semi-technical vocabulary*. CLCS Occasional Paper No. 25 Trinity College.
- Flowerdew, J. 1993. 'Concordancing as a tool in course design.' *System* 21: 231-44.
- Higgins, J. J. 1966. 'Hard facts.' *ELT Journal* 21: 55-60.
- Kövecses, Z. and P. Szabo. 1996. 'Idioms: a view from cognitive semantics.' *Applied Linguistics* 17: 326-55.
- Krejcie, R. and D. Morgan. 1970. 'Determining sample size for research activities.' *Educational and Psychological Measurement* 30: 607-10.
- Lakoff, G. 1987. *Women, Fire and Dangerous Things*. Chicago: University of Chicago Press.
- Martin, A. V. 1976. 'Teaching academic vocabulary to foreign graduate students.' *TESOL Quarterly* 10: 91-7.
- Nagy, W. E. and Anderson, R. C. 1984. 'How many words are there in printed school English?' *Reading Research Quarterly* 19: 304-30.
- Nagy, W. E., R. Anderson, M. Schommer, J. A. Scott, and A. Stallman. 1989. 'Morphological families in the internal lexicon.' *Reading Research Quarterly* 24: 263-82.
- Nation, I. S. P. and K. Hwang. 1995. 'Where would general service vocabulary stop and special purposes vocabulary begin?' *System* 23: 35-41.
- Pearsall, J. (ed.) 1998. *The New Oxford Dictionary of English*. Oxford: Clarendon Press.
- Rosenthal, R. 1987. *Judgment Studies: Design, Analysis and Meta-analysis*. Cambridge: Cambridge University Press.
- Ruhl, C. 1989. *On Monosemy—A Study in Linguistic Semantics*. Albany: State University of New York Press.
- Scott, M. 1999. WordSmith Tools. Version 3.00.00. Oxford: Oxford University Press (<http://www.oup.co.uk/>).
- Sinclair, J. M. 1991. *Corpus, Concordance, Collocation*. Oxford: Oxford University Press.
- Visser, A. 1989. 'Learning core meanings.' *Guidelines* 11/2: 10-17.
- West, M. 1953. *A General Service List of English Words*. London: Longman, Green and Co.
- Wittgenstein, L. 1953. *Philosophical Investigations*. New York: Macmillan.
- Yang H. 1986. 'A new technique for identifying scientific/technical terms and describing science texts.' *Literary and Linguistic Computing* 1: 93-103.