## Corpus design of corpora of living languages:

- Non-exhaustive
- Representativeness is an issue
- Sampling is unavoidable
- Balance and sampling are to be considered to ensure representativeness

## What sort of corpus is the BNC?

**Monolingual**: It deals with modern British English, not other languages used in Britain. However non-British English and foreign language words do occur in the corpus.

**Synchronic**: It covers British English of the late twentieth century, rather than the historical development which produced it.

**General**: It includes many different styles and varieties, and is not limited to any particular subject field, genre or register. In particular, it contains examples of both spoken and written language.

**Sample**: For written sources, samples of 45,000 words are taken from various parts of single-author texts. Shorter texts up to a maximum of 45,000 words, or multi-author texts such as magazines and newspapers, are included in full. Sampling allows for a wider coverage of texts within the 100 million limit, and avoids over-representing idiosyncratic texts.

**BNC-Guide** 

## Unit A2 Representativeness, balance and sampling

#### **A2.1 INTRODUCTION**

We noted in Unit A1 that representativeness is an essential feature of a corpus. It is this feature that is typically used to distinguish a corpus from an archive (i.e. a random collection of texts). A corpus is designed to represent a particular language or language variety whereas an archive is not. Unless you are studying a dead language or highly specialized sub-language (see Unit A2.3 for further discussion), it is virtually impossible to analyse every extant utterance or sentence of a given language. Hence, sampling is unavoidable. Yet how can you be sure that the sample you are studying is representative of the language or language variety under consideration? The answer is that one must consider balance and sampling to ensure representativeness. Hence, this unit introduces the key concept of corpus representativeness as well as the related issues of balance and sampling. We will first explain what we mean by *representativeness* (Unit A2.2), followed by a discussion of the representativeness of general and specialized corpora (Unit A2.3). We will then move on to discuss corpus balance (Unit A2.4) and finally introduce sampling techniques (Unit A2.5).

## Representativeness

What does representativeness s mean in corpus linguistics?

According to Leech (1991:27), a corpus is thought to be representative of the language variety it is supposed to represent if the findings based on its contents can be generalized to the said language variety.

Biber (1993: 243) defines representativeness from the viewpoint of how this quality is achieved: 'Representativeness refers to the extent to which a sample includes the full range of variability in a population.'

A corpus is essentially a *sample* of a language or language variety (i.e. *population*). Sampling is entailed in the compilation of virtually any corpus of a living language. In this respect, the representativeness of most corpora is to a great extent determined by two factors:

- the range of genres included in a corpus (i.e. *balance*, see Unit *A2A*) and how the
- text chunks for each genre are selected (i.e. *sampling*, see Unit A2.5).

## **Ballance**

As with representativeness, the acceptable balance of a corpus is determined by its intended uses.

- Hence, a general corpus which contains both written and spoken data (e.g. the BNC, see Unit A7.2) is balanced;
- so are written corpora such as Brown and LOB (see Unit A7.4),
- and spoken corpora like CANCODE (see Unit A7.5);
- domain-specific corpora (e.g, the HKUST Computer Science Corp us, see Unit A7.3) can also claim to be balanced.

A balanced corpus usually covers a wide range of text categories which are supposed to be representative of the language or language variety under consideration. These text categories are typically sampled proportionally (see UnitA2.5) for inclusion in a corpus so that ' it offers a manageably small scale model of the linguistic material which the corpus builders wish to study' (Atkins et al 1992: 6).

## **Sampling**

In order to obtain a representative sample from a population, the first concern to be addressed is to define the *sampling unit* and the boundaries of the population. For written text, for example, a sampling unit may be a book, periodical or newspaper. The population is the assembly of all sampling units while the list of sampling units is referred to as a *sampling frame*.

In corpus design, a population can be defined in terms of language production, language reception or language as a product. The first two designs are basically demographically oriented as they use the demographic distribution (e.g, age, sex, social class) of the individuals who produce/ receive language data to define the population while the last is organized around text category/genre of language data.

However, it can be notoriously difficult to define a population or construct a sampling frame, particularly for spoken language, for which there are no ready-made sampling frames in the form of catalogues or bibliographies.

## **Sampling**

Once the target population and the sampling frame are defined, different sampling techniques can be applied to choose a sample which is as representative as possible' of the population.

- A basic sampling method is *simple random sampling*.
- [Another method is] *stratified random sampling*.

A further decision to be made in sampling relates to sample size. For example, with written language, should we sample full texts (i.e. whole document s) or text chunks? If text chunks are to be sampled, should we sample text initial, middle or end chunks?

Another sampling issue, which particularly relates to stratified sampling, is the proportion and number of samples for each text category. The numbers of samples across text categories should be proportional to their frequencies and/or weights in the target population in order for the resulting corpus to be considered as representative.

## **A2 Summary**

This unit introduced some important concepts in corpus linguistics - representativeness, balance and sampling. A corpus is considered representative if what we find on the basis of the corpus also holds for the language or language variety it is supposed to represent. For most corpora, representativeness is typically achieved by balancing, i.e. covering a wide variety of frequent and important text categories that are proportion ally sampled from the target population. Claims of corpus representativeness and balance, however, should be interpreted in relative terms and considered as a statement of faith rather than as fact, as presently there is no objective way to balance a corpus or to measure its representativeness. Furthermore, it is only by considering the research question one has to address that one is able to determine what is an acceptable balance for the corpus one should use and whether it is suitably representative. The concepts introduced in this unit will help you to determine if a particular corpus is suitable for your intended research. They are also helpful in determining whether a research question is amenable to corpus analysis.

There is a broad consensus among the participants in the project and among corpus linguists that a general-purpose corpus of the English language would ideally contain a high proportion of spoken language in relation to written texts. However, it is significantly more expensive to record and transcribe natural speech than to acquire written text in computer-readable form. Consequently the spoken component of the BNC constitutes approximately 10 per cent (10 million words) of the total and the written component 90 per cent (90 million words). These were agreed to be realistic targets, given the constraints of time and budget, yet large enough to yield valuable empirical statistical data about spoken English. In the BNC sampler, a two per cent sample taken from the whole of the BNC, spoken and written language are present in approximately equal proportions, but other criteria are not equally balanced.

From the start, a decision was taken to select material for inclusion in the corpus according to an overt methodology, with specific target quantities of clearly defined types of language. This approach makes it possible for other researchers and corpus compilers to review, emulate or adapt concrete design goals. This section outlines these design considerations, and reports on the final make-up of the BNC.

This and the other tables in this section show the actual make-up of the second version of the British National Corpus (the BNC World Edition) in terms of

- texts: number of distinct samples not exceeding 45,000 words
- S-units: number of <s> elements identified by the CLAWS system (more or less equivalent to sentences)
- W-units: number of <w> elements identified by the CLAWS system (more or less equivalent to words)

The XML Edition of the BNC contains 4049 texts and occupies (including all markup) 5,228,040 Kb, or about 5.2 Gb. In total, it comprises just under 100 million orthographic words (specifically, 96986707), but the number of w-units (POS-tagged items) is slightly higher at 98363783. The tagging distinguishes a further 13614425 punctuation strings, giving a total content count of 110691482 strings. The total number of s-units tagged is about 6 million (6026284). Counts for these and all the other elements tagged in the corpus are provided in the corpus header.

In the following tables both an absolute count and a percentage are given for all the counts. The percentage is calculated with reference to the relevant portion of the corpus, for example, in the table for "written text domain", with reference to the total number of w-units in written texts. Note that punctuation strings are not included in these totals. The reference totals used are given in the first table below.

## Table 1. Text type

	<b>~</b> ,				
	texts	w-units	%	s-units	%
Spoken demographic	153	4233955	4.30	610557	10.13
Spoken context-governed	<b>755</b>	6175896	6.27	427523	7.09
Written books and periodicals	2685	79238146	80.55	4395581	72.94
Written-to-be-spoken	35	1278618	1.29	104665	1.73
Written miscellaneous	421	7437168	7.56	487958	8.09

All texts are also classified according to their date of production. For spoken texts, the date was that of the recording. For written texts, the date used for classification was the date of publication of the source edition used, for the most part; in the case of imaginative works, however, the date of first publication of the work was used. Informative texts were selected only from 1975 onwards, imaginative ones from 1960, reflecting their longer 'shelf-life', though most (75 per cent ) of the latter were published no earlier than 1975.

### **Table 2. Publication date**

texts w-units % s-units %
Unknown 162 1831585 1.86 126416 2.09
1960-197446 1718449 1.74 119510 1.98
1975-1984169 4730889 4.80 257962 4.28
1985-199336729008286091.58552239691.63

(http://www.natcorp.ox.ac.uk/XMLedition/URG/BNCdes.html)

Sampling basis: production and reception

While it is sometimes useful to distinguish in theory between language which is received (read and heard) and that which is produced (written and spoken), it was agreed that the selection of samples for a general-purpose corpus must take account of both perspectives.

Text that is published in the form of books, magazines, etc., is not representative of the totality of written language that is produced, as writing for publication is a comparatively specialized activity in which few people engage.

However, it is much more representative of written language that is received, and is also easier to obtain in useful quantities, and thus forms the greater part of the written component of the corpus.

There was no single source of information about published material that could provide a satisfactory basis for a sampling frame, but a combination of various sources furnished useful information about the totality of written text produced and, particularly, received, some sources being more significant than others.

They are principally statistics about books and periodicals that are published, bought or borrowed.

Catalogues of books published per annum tell us something about production but little about reception as many books are published but hardly read.

A list of books in print provides somewhat more information about reception as time will weed out the books that nobody bought (or read): such a list will contain a higher proportion of books that have continued to find a readership.

The books that have the widest reception are presumably those that figure in **bestseller lists**, particularly **prize winners** of competitions such as the Booker or Whitbread.

Such works were certainly candidates for inclusion in the corpus, but the statistics of book-buying are such that very few texts achieve high sales while a vast number sell only a few or in modest numbers.

If texts had been selected in strict arithmetical proportion to their sales, their range would have been severely limited. However, where a text from one particular subject domain was required, it was appropriate to prefer a book which had achieved high sales to one which had not.

Library lending statistics, where these are available, also indicate which books enjoy a wide reception and, like lists of books in print, show which books *continue* to be read.

Similar observations hold for magazines and periodicals. lists of current magazines and periodicals are similar to catalogues of published books, but perhaps more informative about language reception, as it may be that periodicals are bought and read by a wider cross-section of the community than books. Also, a periodical that fails to find a readership will not continue to be published for long.

Periodical circulation figures have to be treated with the same caution as bestseller lists, as a few titles dominate the market with a very high circulation. To concentrate too exclusively on these would reduce the range of text types in the corpus and make contrastive analysis difficult.

Published written texts were selected partly at random from Whitaker's *Books in Print* for 1992 and partly systematically, according to the selection features outlined in section Selection features below.

Available sources are concerned almost exclusively with published books and periodicals. It is much more difficult to obtain data concerning the production or reception of unpublished writing. Intuitive estimates were therefore made in order to establish some guidelines for text sampling in the latter area.

#### **Selection features**

Texts were chosen for inclusion according to three selection features:

- domain (subject field),
- time (within certain dates) and
- medium (book, periodical, etc.).

The purpose of these selection features was to ensure that the corpus contained a broad range of different language styles, for two reasons. The first was so that the corpus could be regarded as a microcosm of current British English in its entirety, not just of particular types. The second was so that different types of text could be compared and contrasted with each other.

#### **Selection Procedure**

Each selection feature was divided into classes (e.g. 'Medium' into books, periodicals, unpublished etc.; 'Domain' into imaginative, informative, etc.) and target percentages were set for each class.

These percentages are quite independent of each other: there was no attempt, for example, to make 25 per cent of the selected periodicals imaginative.

The design proposed that seventy-five per cent of the samples be drawn from informative texts, and the remaining 25 per cent from imaginative texts.

It further proposed that titles be taken from a variety of media, in the following proportions:

- 60 per cent from books,
- 30 per cent from periodicals,
- 10 per cent from miscellaneous sources (published, unpublished, and written to be spoken).

Half of the books in the 'Books and Periodicals' class were selected at random from Whitaker's *Books in Print 1992*. This was to provide a control group to validate the categories used in the other method of selection: the random selection disregarded Domain and Time, but texts selected by this method were classified according to these other features after selection.

Sample size and method For books, a target sample size of 40,000 words was chosen.

No extract included in the corpus exceeds 45,000 words.

For the most part, texts which in their entirety were shorter than 40,000 words were further reduced by ten per cent for copyright reasons; a few texts longer than the target size were however included in their entirety.

Text samples normally consist of a continuous stretch of discourse from within the whole. A convenient breakpoint (e.g. the end of a section or chapter) was chosen as far as possible to begin and end the sample so that high-level discourse units were not fragmented.

Where possible, no more than one sample was taken from any one text; for newspaper texts and large encyclopaedic works, no sample greater than 40,000 words was taken.

Samples were taken randomly from the beginning, middle or end of longer texts. (In cases where a publication included essays or articles by a variety of authors of different nationalities, the work of non-UK authors was omitted.)

Some types of written material are composite in structure: that is, the physical object in written form is composed of more than one text unit. Important examples are issues of a newspaper or magazine which, though editorially shaped as a document, contain discrete texts, each with its specific authorship, stylistic characteristics, register and domain.

The BNC attempts to separate these discrete texts where appropriate and to classify them individually according to the selection and classification features. As far as possible, the individual stories in one issue of a newspaper were grouped according to domain, for example as 'Business' articles, 'Leisure' articles, etc.

The following subsections discuss each selection criterion, and indicate the actual numbers of words in each category included.

#### **Domain**

Classification according to subject field seems hardly appropriate to texts which are fictional or which are generally perceived to be literary or creative.

Consequently, these texts are all labelled imaginative and are not assigned to particular subject areas.

All other texts are treated as informative and are assigned to one of the eight domains listed below.

(http://www.natcorp.ox.ac.uk/XMLedition/URG/BNCdes.html)

**Table 3. Written Domain** 

	texts	w-units	%	s-units	%
Imaginative	476	16496420	18.75	1352150	27.10
Informative: natural & pure science	146	3821902	4.34	183384	3.67
Informative: applied science	370	7174152	8.15	356662	7.15
Informative: social science	526	14025537	15.94	698218	13.99
Informative: world affairs	483	17244534	19.60	798503	16.00
Informative: commerce & finance	295	7341163	8.34	382374	7.66
Informative: arts	261	6574857	7.47	321140	6.43
Informative: belief & thought	146	3037533	3.45	151283	3.03
Informative: leisure	438	12237834	13.91	744490	14.92

The evidence from catalogues of books and periodicals suggests that imaginative texts account for significantly less than 25 per cent of published output, and unpublished reports, correspondence, reference works and so on would seem to add further to the bulk of informative text which is produced and consumed. However, the overall distribution between informative and imaginative text samples is set to reflect the influential cultural role of literature and creative writing. The target percentages for the eight informative domains were arrived at by consensus within the project, based loosely upon the pattern of book publishing in the UK during the past 20 years or so, as reflected in the categorized figures for new publications that appear annually in Whitaker's *Book list*.

#### Medium

This categorisation is broad, since a detailed taxonomy or feature classification of text medium could have led to such a proliferation of subcategories as to make it impossible for the BNC adequately to represent all of them.

The labels used here are intended to be comprehensive in the sense that any text can be assigned with reasonable confidence to these macro categories.

The labels we have adopted represent the highest levels of a fuller taxonomy of text medium.

**Table 4. Written Medium** 

	texts	w-units	%	s-units	%
Book	1411	50293803	57.18	2887523	57.88
Periodical	1208	28609494	32.52	1487644	29.82
Miscellaneous published	238	4233135	4.81	287700	5.76
Miscellaneous unpublished	249	3538882	4.02	220672	4.42
To-be-spoken	35	1278618	1.45	104665	2.09

The 'Miscellaneous published' category includes brochures, leaflets, manuals, advertisements. The 'Miscellaneous unpublished' category includes letters, memos, reports, minutes, essays. The 'written-to-bespoken' category includes scripted television material, play scripts etc.

Descriptive features
Written texts may be further classified according to sets of descriptive features.

These features *describe* the sample texts; they did not determine their selection.

This information is recorded to allow more delicate contrastive analysis of particular sets of texts.

As a simple example, the gross division into two time periods in the selection features can, of course, be refined and subcorpora defined over the BNC for more specific dates.

However, the relative sizes of such subcorpora are undefined by the BNC design specification.

These descriptive features were monitored during the course of the data gathering, and text selection, in cases where a free choice of texts was available, took account of the relative balance of these features.

Thus although no relative proportions were defined for different target age groups (for example), we ensured that the corpus does contain texts intended for children as well as for adults.

The following tables summarize the results for the first release of the corpus. Note that many texts remain unclassified.

#### **Author information**

Information about authors of written texts was included only where it was readily available, for example from the dust-wrapper of a book.

Consequently, the coverage of such information is very patchy.

The authorship of a written text was characterized as 'corporate' where it was produced by an organization and no specific author was given, and as 'multiple' in cases where several authors were named.

Author sex was classified as 'mixed' where more than one author of either sex was specified, and 'unknown' where it could not reliably be determined from the author's name.

Note that 'author age' means the author's age at the time of creation of the work concerned.

(http://www.natcorp.ox.ac.uk/XMLedition/URG/BNCdes.html)

## Table 5. Author type

texts w-units % s-units %
Unknown 211 3786835 4.30 174371 3.49
Corporate author 347 6497144 7.38 455649 9.13
Multiple author 1322 34563219 39.29 1810901 36.30

Sole author 1261 43106734 49.01 2547283 51.06

(http://www.natcorp.ox.ac.uk/XMLedition/URG/BNCdes.html)

#### Table 6. Sex of author

texts w-units % s-units %

Unknown 1573 36161115 41.11 1968162 39.45

Author sex Male 920 30665582 34.86 1671420 33.50

Author sex Female 414 14588260 16.58 967522 19.39

Author sex Mixed 234 6538975 7.43 381100 7.64

(http://www.natcorp.ox.ac.uk/XMLedition/URG/BNCdes.html)

Table 7. Author age-group

	texts	w-units	%	s-units	%
Unknown	2518	66000719	75.04	3687586	73.92
Author age 0-14	3	59559	0.06	3443	0.06
Author age 15-24	19	542578	0.61	29810	0.59
Author age 25-34	66	2267139	2.57	159455	3.19
Author age 35-44	191	6726926	7.64	410143	8.22
Author age 45-59	205	7230714	8.22	410644	8.23
Author age 60+	139	5126297	5.82	287123	5.75

Table 8. Domicile

	texts	w-units	%	s-units	%
Unknown	2272	57227155	65.06	3133068	62.80
Author domicile UK and Ireland	841	29760000	33.83	1798301	36.05
<b>Author domicile Commonwealth</b>	12	411207	0.46	25759	0.51
<b>Author domicile Continental Europe</b>	6	234402	0.26	12466	0.24
Author domicile USA	8	245604	0.27	15675	0.31
Author domicile Elsewhere	2	75564	0.08	2935	0.05

# Target audience

Some attempt was made to characterize the kind of audience for which written texts were produced in terms of

- age,
- sex and
- 'level' (a subjective assessment of the text's technicality or difficulty).

The last of these proved very difficult to assess and was very frequently confused with circulation size or audience size; for that reason, no figures for it are included here.

Table 9. Audience age

	texts	w-units	%	s-units	%
Child audience	42	903690	1.02	81074	1.62
Teenager audience	78	1831178	2.08	138098	2.76
Adult audience	2911	81928776	93.14	4597388	92.16
Any audience	110	3290288	3.74	171644	3.44

(http://www.natcorp.ox.ac.uk/XMLedition/URG/BNCdes.html)

### Table 10. Audience sex

texts w-units % s-units %
Unknown 706 20271270 23.04 1131254 22.67
Male audience 61 2396935 2.72 135950 2.72
Female audience 175 6904137 7.84 503629 10.09
Mixed audience 2199 58381590 66.37 3217371 64.49

(http://www.natcorp.ox.ac.uk/XMLedition/URG/BNCdes.html)

Miscellaneous classification information: Written texts were also characterized according to their place of publication and the type of sampling used.

**Table 11. Publication place** 

	texts	w-units	%	s-units	%
Unknown	690	14718827	16.73	788440	15.80
UK (unspecific) publication	263	7163111	8.14	380824	7.63
Ireland publication	37	570652	0.64	31793	0.63
UK: North (north of Mersey-Humber line) publication	191	3781055	4.29	228247	4.57
UK: Midlands (north of Bristol Channel-Wash line) publication	93	2590345	2.94	177308	3.55
UK: South (south of Bristol Channel-Wash line) publication	1853	58587808	66.61	3360401	67.36
United States publication	14	542134	0.61	21191	0.42

(http://www.natcorp.ox.ac.uk/XMLedition/URG/BNCdes.html)

# **Table 12. Sampling type**

	texts	w-units	%	s-units	%
Unknown	1583	3555110	0240.4	2 1991798	39.93
Whole text	270	652497	5 7.41	433722	8.69
<b>Beginning sample</b>	584	2107522	22 23.9	61119251	22.43
Middle sample	510	1845480	0720.9	8 1049692	21.04
End sample	119	4317320	6 4.90	253322	5.07
Composite sample	75	2030500	0 2.30	140419	2.81

In addition to the above, standard bibliographic details such as author, title, publication details, extent, topic keywords etc. were recorded for the majority of texts, as further described below (see 5 The header).

Selection procedures employed Books

Roughly half the titles were randomly selected from available candidates identified in Whitaker's *Books in Print* (BIP), 1992, by students of Library and Information Studies at Leeds City University. Each text randomly chosen was accepted only if it fulfilled certain criteria: it had to be published by a British publisher, contain sufficient pages of text to make its incorporation worthwhile, consist mainly of written text, fall within the designated time limits, and cost less than a set price. The students noted the ISBN, author, title and price of each book thus selected; the final selection weeded out texts by non-UK authors.

Half of the books having been selected by this method, the remaining half were selected systematically to make up the target percentages in each category. The selection proceeded as follows.

(http://www.natcorp.ox.ac.uk/XMLedition/URG/BNCdes.html)

## **Bestsellers**

Because of their wide reception, bestsellers were obvious candidates for selection. The lists used were those that appeared in the *Bookseller* at the end of the years 1987 to 1993 inclusive. Some of the books in the lists were rejected, for a variety of reasons. Obviously books that had already been selected by the random method were excluded, as were those by non-UK authors. In addition, a limit of 120,000 words from any one author was imposed, and books belonging to a domain or category whose quota had already been reached were not selected. Other bestseller lists were obtained from *The Guardian*, the British Council, and from Blackwells Paperback Shop.

The titles yielded by this search were mostly in the Imaginative category.

# **Literary prizes**

The criteria for inclusion were the same as for bestsellers. The prize winners, together with runners-up and shortlisted titles, were taken from several sources, principally Anne Strachan, *Prizewinning literature: UK literary award winners,* London, 1989. For 1990 onwards the sources used were: the last issue of the *Bookseller* for each year; *The Guardian Index, 1989*–, entries under the term 'Literature'; and *The Times Index, 1989-*, entries under the term 'Literature — Awards'.

Literary prizes are in the main awarded to works that fall into the Imaginative category, but there are some Informative ones also.

(http://www.natcorp.ox.ac.uk/XMLedition/URG/BNCdes.html)

## **Library loans**

The source of statistics in this category was the record of loans under Public Lending Right, kindly provided by Dr J. Parker, the Registrar. The information comprised lists of the hundred most issued books and the hundred most issued children's books, in both cases for the years 1987 to 1993.

The lists consist almost exclusively of imaginative literature, and many titles found there also appear in the lists of bestsellers and prize winners.

## **Additional texts**

As collection proceeded, monitoring disclosed potential shortfalls in certain domains. A further selection was therefore made, based on the 'Short Loan' collections of seven University libraries. (Short Loan collections typically contain books required for academic courses, which are consequently in heavy demand.)

## Periodicals and magazines

Periodicals, magazines and newspapers account for 30 per cent of the total text in the corpus. Of these, about 250 titles were issues of newspapers. These were selected to cover as wide a spectrum of interests and language as possible. Newspapers were selected to represent as wide a geographic spread as possible: *The Scotsman* and the *Belfast Telegraph* are both represented, for example.

### Other media

In addition to samples from books, periodicals, and magazines, the written part of the corpus contains about seven million words classified as 'Miscellaneous Published', 'Miscellaneous Unpublished', or as 'Written to be spoken'. The distinction between 'published' and 'unpublished' is not an easy one; the former category largely contains publicity leaflets, brochures, fact sheets, and similar items, while the latter has a substantial proportion of school and university essays, unpublished creative writing or letters, and internal company memoranda. The 'written to be spoken' material includes scripted material, intended to be read aloud such as television news broadcasts; transcripts of more informal broadcast materials such as discussions or phone-ins are included in the spoken part of the corpus.

# **Copyright permissions**

Before a selected text could be included, permissions had to be obtained from the copyright owner (publisher, agent, or author). A standard Permissions Request was drafted with considerable care, but some requests were refused, or simply not answered even after prompting, so that the texts concerned had to be excluded or replaced.