

Corpus Design Criteria

Sue Atkins
Jeremy Clear
Nicholas Ostler

15th January 1991

Contents

Introduction	1
1 Defining Text Collections and a Unit of Text	1
1.1 Texts	2
1.1.1 Written Texts	2
1.1.2 Spoken Texts	3
2 Stages in Corpus Building: Tasks, Expertise, Personnel	3
2.1 Planning	3
2.2 Permissions	4
2.3 Data Capture	4
2.4 Text Handling	4
2.4.1 Basic Tools	5
2.4.2 Advanced Text Handling	5
2.5 User Feedback and Corpus Development	5
3 Corpora and Copyright	6
4 Population and Sampling	6
4.1 Defining the Population	7
4.2 Describing the Population	8
5 Markup	8
5.1 Methodological Considerations	8
5.1.1 Converting Written Material	8
5.1.2 Transcription of Speech	9
5.2 Features for Markup: written text	10
5.3 Features for Markup: spoken text	11
6 Corpus Typology	13
7 Text Typology	14
7.1 Text Attributes	15
7.2 Reference Coding	19
7.3 Case Study: Designing a Lexicographic Corpus	20
7.3.1 A Taxonomy of Text Types	20
7.3.2 Topic	21

7.3.3	Controlled Parameters for Written Texts	22
8	Progress to date with Standards for Corpora	23
8.1	Terminology	23
8.2	Encoding Standards	23
8.3	Evaluation Standards	25
9	Potential Users and Uses	25
9.1	Generalities	25
9.2	Language Specialists	26
9.3	Content Specialists	27
9.4	Media Specialists	28
9.5	Types of Use	28
	Select Bibliography	30

Introduction

There has been over the past few years a tremendous growth in interest and activity in the area of corpus building and analysis. European, USA and Japanese efforts in the development of NLP and IT are converging on the recognition of the importance of some sort of corpus-based research as part of the infrastructure for the development of advanced language processing applications. Statistical processing of text corpora has been demonstrated as a viable approach to some of the traditional hard problems of computational linguistics, machine translation and knowledge engineering.

Our aim in this paper is to identify the principal aspects of corpus creation and the major decisions to be made when creating an electronic text corpus, for whatever purpose; and to discuss in detail the criteria that must inform some of these decisions and that are relevant to the establishment of basic standards for corpus design. We expect the paper to form the basis of an in-depth study of corpus design criteria, with the object of defining the minimal set necessary to foster the creation of high-quality compatible corpora of different languages, for different purposes, in different locations, and using different types of software and hardware.

To this end, we attempt to identify the principal features in corpus design, and to note others which must not be forgotten but which need not be addressed in the initial stages of corpus building. Our aim is not to make a comprehensive and water-tight listing of everything it is possible to decide, for we believe that this would be totally counterproductive: rather, a corpus may be thought of as organic, and must be allowed to grow and live if it is to reflect a growing, living language. Our aim is simply to pick out the principal decision points, and to describe the options open to the corpus-builder at each of these points.

We specifically exclude from this paper all consideration of acoustic corpora, which falls outside our brief; however this is not to imply that research into spoken language is any less important; simply that the needs of the two groups are different, making it possible to consider these two aspects of NLP separately. What we say in this paper relates to the needs of speech research only insofar as any general text corpus must include as high as possible proportion of transcribed spoken text, which is included in our understanding of the terms ‘text’ or ‘lexical material’.

1 Defining Text Collections and a Unit of Text

We distinguish four types of text collection¹, which we find helpful and urge the community to accept:

archive a repository of readable electronic texts not linked in any coordinated way.

electronic text library (or ETL, Fr. ‘textothèque’) a collection of electronic texts in standardised format with certain conventions relating to content etc, but without rigorous selectional constraints.

corpus a subset of an ETL, built according to explicit design criteria for a specific purpose, eg the Corpus Révolutionnaire (Bibliothèque Beaubourg, Paris), the Cobuild Corpus, the Longman/Lancaster corpus, the Oxford Pilot corpus.

subcorpus a subset of a corpus, either a static component of a complex corpus or a dynamic selection from a corpus during on-line analysis.

This paper is concerned with ETLs, corpora and subcorpora but for the sake of brevity we use the word corpus to refer to all three types of collection.

¹The first three types are those identified by Bernard Quemada

1.1 Texts

A corpus which is designed to constitute a representative sample of a defined language type will be concerned with the sampling of *texts*. For the purposes of studying spoken language in transcription (not speech *per se*) it is convenient to use the term ‘text’ to include transcribed speech. The use of the word to describe a unit of text, informally considered to be integral in some way, raises some issues of definition for the corpus builder. If, for example, it is decided that the corpus should include ‘full texts’, or if sampling is to be carried out by some random selection of texts, then it will be important to consider what is meant by a text. Since the notion of text is derived strongly from models of written language, and is more tenuously applied to speech events, we should consider the definition separately for these two modes.

1.1.1 Written Texts

The printed monograph or work of narrative fiction is the model for the notion of a ‘text’. It manifests several apparently criterial characteristics:

- it is discursive and typically at least several pages long;
- it is integral;
- it is the conscious product of a unified authorial effort;
- it is stylistically homogeneous.

‘Texts’ are often assumed to be a series of coherent sentences and paragraphs. By integral, we mean that a printed book usually has a beginning, middle and end and is considered to be complete in itself. Even though a book may have more than one author, their collaboration usually counts as authorial unity (no particular parts of the book are ascribed to any one individual) and the result is usually a unit of stylistically consistent language.

The important aspects in these criterial features for the corpus builder are those relating to stylistics and text linguistics. The designer of the corpus will wish to neutralise as far as possible the effects of sampling bias and the stylistic idiosyncrasies of one particular author can be reduced in significance if texts by many different authors are included. The need to control stylistic parameters leads to the concern with a unified authorial effort and consistent style. Similarly, if the corpus is to provide the basis for studies of cohesion, discourse analysis and text linguistics—all linguistic patterning beyond the sentence or paragraph—then the integrity of the samples as textual units ought to be taken into consideration.

Novels fit the above-mentioned schema very neatly and are prototypical ‘texts’, but there are many types of written language which are more problematic. Listed below are some examples of language units which are likely to be incorporated in a general corpus and which illustrate typical deviations from the model instance.

Small ads in newspapers where the corpus builder might prefer to make a *collection* of these small ads and treat this as one text.

An article in a newspaper or magazine: it may be convenient to treat one issue of a newspaper (and single issues of other periodicals) as one text.

An academic *Festschrift*, learned journal, etc. where the bibliographic data applies to the whole book but the papers differ linguistically to such an extent that they might be best treated as a discrete text.

A poem It is often more convenient to gather many short poems by the same author into collections and to treat each collection as a text.

Published correspondence in which the letters are the product of two authors, but they constitute a single discourse. The critical apparatus, introduction and editorial material will be yet another author's intervention. The way this problem is handled will depend on the requirements of the corpus.

1.1.2 Spoken Texts

The difficulty and high cost of recording and transcribing natural speech events leads the corpus linguist to adopt a more open strategy to collecting spoken language. It may be more acceptable to the corpus builder to capture *fragments* of speech than fragment of writing. A stretch of speech can be thought of as forming a text, rather than a fragment, if one the following conditions applies:

- the speech unit starts when the participants come together and ends when they part;
- the speech has an obvious opening and closing;

Some examples of units of speech which might be considered to be texts are:

- an informal face-to-face conversation**
- a telephone conversation**
- a lecture**
- a meeting**
- an interview**
- a debate.**

2 Stages in Corpus Building: Tasks, Expertise, Personnel

2.1 Planning

The planning stage of corpus building will aim to arrive at specifications in two related areas: the linguistic design of the corpus and the project costs and administration.

The linguistic design will need to establish at the least what type of corpus is being constructed (see Section 6, Corpus Typology), the sizes of the text samples to be included, the range of language varieties (synchronic) and the time period (diachronic) to be sampled, whether to include writing and speech and the approximate level of encoding detail to be recorded in electronic form.

The administrative planning will be concerned with the costs and the stages of the corpus building work. The primary stages are

- Specifications and design
- Selection of sources
- Obtaining copyright permissions
- Data capture and encoding/markup
- Corpus processing

The design will involve consultation with representatives of the anticipated users of the corpus and also perhaps with other specialists. For a corpus which aims to represent general language synchronically, a sociolinguist may be called upon. Decisions concerning the sampling strategy may involve a statistician and statistical expertise will almost certainly be valuable to some extent throughout the corpus building project. Linguistic expertise will also be required to ensure that all decisions concerning the design, balance, encoding and processing of the corpus are appropriate to the linguistic aims of the corpus project. It is because of the particular *linguistic* interest in a large body of computerised text that a

language corpus is quite different from any of the very large on-line information databases which are available (commercially and for research purposes) around the world.

Estimates will be required of the hardware and software needs of the corpus project. These will depend to a large extent on the amount of processing and manipulation of the corpus which is to be carried out.

2.2 Permissions

The selection of sources might be based on a systematic analysis of the target population or on a random selection method. Alternatively, the need for large volumes of data may lead one to adopt a more opportunistic approach to the collection of text. In most cases, copyright permission will need to be obtained for texts to be computerised. This can be a time-consuming and increasingly difficult operation. Co-ordinated effort is required to ensure that the texts are used in accordance with copyright legislation (see Section 3, Corpora and Copyright). It is wise to seek legal advice in this area, even when the data being stored on computer is apparently out of copyright, or is to be used only by university researchers.

2.3 Data Capture

Data capture will also be time-consuming and the costs incurred will be unavoidable and determined to a large extent by the amount of text to be captured. Printed material can be scanned by OCR devices with appropriate software. This method of conversion from print to electronic form is becoming more cost effective each year as OCR technology improves. Alternatively, text can be keyboarded manually. For transcription of audio recordings and for the capture of degraded or complex printed material keyboarding is the only option.

Text may also be acquired in machine-readable form (either having been captured by someone else or having originated on computer). This eliminates the costs of data capture but may require time-consuming parsing, reformatting and restructuring to bring the data into line with the corpus conventions for encoding and markup.

The texts in the corpus will probably (though not necessarily) be marked up with embedded codes to signal elements of structure and to record features of the original text source. This is an important issue and one which will have implications for the scheduling and costing of a corpus project (see Section 5, Markup). There is a growing consensus that SGML provides a suitable basis for a standard markup scheme for texts held of computer and the Text Encoding Initiative (TEI) published in July 1990 its draft guidelines for the encoding and interchange of machine-readable texts.

Whatever method of data capture is adopted, the text will require some degree of validation and error-correction to ensure that it is reasonably accurate and consistent with the encoding conventions for the corpus.

2.4 Text Handling

The mere existence of a large corpus will not satisfy demand for linguistic data. A set of general tools for processing the corpus will be essential. Many such tools already exist and are in use, but they are often designed to meet very specific local needs and there is work to be done on agreeing on standard formats for data derived from a corpus, for word-class tagging, parse tree notations, semantic labelling, etc. KWIC concordances, word frequency lists, collocation statistics will be basic requirements and software to perform these basic tasks can be made widely available as part of the corpus package. Tagging software and

parsers will be more difficult to implement and their design more contentious, but these tools should also be available.

2.4.1 Basic Tools

These might include:

Word Frequency Software to produce lists of word types and their frequency in the corpus and perhaps also some statistical profile of the relation of types to tokens in the corpus, indications of the distribution of types across the text categories, and graphical displays to summarise these lists in a form which can be assimilated by the user of the corpus.

Concordancing Text retrieval and indexing software (of the kind already provided by such packages as Word Cruncher, PAT, TACT, Free Text Browser, SGML Search, etc.) with features appropriate for linguistic analysis.

2.4.2 Advanced Text Handling

In order to allow more sophisticated statistical analyses to be carried out on a large corpus, it may be useful to implement a number of more advanced text processing tools which can automatically process linguistic information in a corpus. Such software might include:

Lemmatisation To relate a particular inflected form of a word to its base form or lemma, thereby enabling the production of frequency and distribution figures which are less sensitive to the incidence of surface strings.

Part-of-speech labelling (sometimes called ‘tagging’) to assign a word class or part-of-speech label to every word. This allows simple syntactic searches to be carried out.

Parsing to assign a fully labelled syntactic tree or bracketing of constituents to sentences of the corpus.

Collocation to compute the statistical association of word forms in the text.

Sense disambiguation (‘homograph separation’) to distinguish which of several possible senses of a given word is being employed at each occurrence. Early work in this area shows promising results by methods based on lookup in a machine-readable dictionary or on identifying collocational patterns.

Link to lexical database to integrate the instances of words or phrases in a corpus with a structured lexical database which records detailed morphological, syntactic, lexical and semantic information about words.

2.5 User Feedback and Corpus Development

Since it is theoretically suspect to aim at achieving a perfectly ‘balanced’ corpus, it is practical to adopt a method of successive approximations. First, the corpus builder attempts to create a representative corpus. Then this corpus is used and analysed and its strengths and weaknesses identified and reported. In the light of this experience and feedback the corpus is enhanced by the addition or deletion of material and the cycle is repeated continually.

There is a need therefore for an appropriate mechanism to allow users of the corpus to relay their findings, comments, warnings and so on back to the corpus development team. This might be done electronically over a network or via e-mail or through regular, scheduled

meetings and discussions between the users and those responsible for the maintenance and development of the corpus. Statistical expertise is likely to be valuable in assessing and testing methods for improving the balance of the corpus to suit the needs of the users.

3 Corpora and Copyright

One of the serious constraints on the development of large text corpora and their widespread use is national and international copyright legislation. It is necessary and sensible to protect, through copyright laws, the rights of authors and publishers in texts that they create. In response to the rapid development of computer technology in the areas of networking, DTP, personal computing and electronic publishing, copyright legislation is being extended and revised to cover what are perceived to be new threats, as well as opportunities, for the writing and publishing industries. The effect for the corpus builder is that it is quite likely that any text (or sample of text) which is to be computerised and included in a corpus will be under copyright protection and that permission will have to be obtained for its use. The following considerations are relevant to copyright and corpora:

- Is the text protected by copyright? National legislation varies but in general texts can pass out of copyright after a certain time period. However, in the US and UK copyright may reside in a particular *edition* of a text (the arrangement and typesetting) even though the original text itself is out of copyright. The computerisation of speech transcriptions may also require permission—particularly if the speech is recorded from radio or TV broadcasts.
- Will payments be offered? Within the publishing world copyright permissions are usually obtained on payment of a fee, a royalty or some combination. It is clear that payment of even modest fees will make the compilation of a corpus of contemporary texts from a large number of different sources so expensive that only a very few organisations will be able to justify the costs. If no fee is to be paid then the copyright holders must be assured that the compilation of a language corpus is no threat whatsoever to the revenue-earning potential of the text and that no direct commercial exploitation is to be made of the corpus.
- What use is to be made of the corpus? If the corpus is to be used for commercial purposes then these will probably need to be clearly stated and defined to the copyright holder. Similarly, if the corpus builder plans to copy and distribute the corpus to other people or to make it accessible by others, this will need careful agreement with the copyright holder.
- How many different sources are to be collected? The use of a renewable, single source, such as the AP Wire Service, requires permission from only one source with perhaps only one fee. If the corpus is to contain many hundreds or thousands of text samples, the administrative and clerical task of identifying copyright holders and obtaining permission can be very substantial.

4 Population and Sampling

In building a natural language corpus one would like ideally to adhere to the theoretical principles of statistic sampling and inference². Unfortunately the standard approaches to

²For a more detailed discussion of this topic see Clear (forthcoming) ‘Corpus sampling’, Proceedings of 11th ICAME Conference, Berlin 1990.

statistical sampling are hardly applicable to building a language corpus. First, often it is very difficult to delimit the total population in any rigorous way. Textbooks on statistical methods almost always focus on clearly defined populations. Second, there is no obvious unit of language which is to be sampled and which can be used to define the population. We may sample words or sentences or ‘texts’ among other things. Third, because of the sheer size of the population and given current and foreseeable resources, it will always be possible to demonstrate that some feature of the population is not adequately represented in the sample. Despite these difficulties, some practical basis for progress can be established. An approach suggested by Woods, Fletcher and Hughes³ is to accept the results of each study as though any sampling had been carried out in the theoretically ‘correct’ way, to attempt to foresee possible objections. In corpus linguistics such a pragmatic approach seems the only course of action. Moreover, there is a tendency to overstate the possibility and effects of experimental error: indeed, good scientific estimation of the possibility and scale of experimental error in statistics of natural language corpora is seldom carried out at all.

All samples are *biased* in some way. Indeed the sampling problem is precisely that a corpus is inevitably biased in some respects. The corpus users must continually evaluate the results drawn from their studies and should be encouraged to report them (see Section 2.5).

The difficulty of drawing firm conclusions when the number of observed instances is few underlines the methodological point made by Wood, Fletcher and Hughes: that researchers should question how the sample was obtained and assess whether this is likely to have a bearing on the validity of the conclusions reached.

4.1 Defining the Population

When a corpus is being set up as a sample with the intention that observation of the sample will allow us to make generalizations about language, then the relationship between the sample and the target population is very important. The more highly specialised the language to be sampled in the corpus, the fewer will be the problems in defining the texts to be sampled. For a general-language corpus however there is a primary decision to be made about whether to sample the language that people hear and read (their *reception*) or the language that they speak and write (their *production*).

Defining the population in terms of language reception assigns tremendous weight to a tiny proportion of the writers and speakers whose language output is received by a very wide audience through the media. However, most linguists would reject the suggestion that the language of the daily tabloid newspapers (though they may have a very wide reception) can be taken to represent the language production of any individual member of the speech community.

The corpus builder has to remain aware of the reception and production aspects, and though texts which have a wide reception are by definition easier to come by, if the corpus is to be a true reflection of native speaker usage, then every effort must be made to include as much production material as possible. For a large proportion of the language community, writing (certainly any extended composition) is a rare language activity. Judged on either of these scales, private conversation merits inclusion as a significant component of a representative general language corpus. Judged in terms of production, personal and business correspondence and other informal written communications form a valuable contribution to the corpus.

³in *Statistics in Language Studies* (1986)

To summarise, we can define the language to be sampled in terms of language production (many producers each with few receivers) and language reception (few producers but each with many receivers). Production is likely to be greatly influenced by reception, but technically only production defines the language variety under investigation. Collection of a representative sample of total language production is not feasible, however. The compiler of a general language corpus will have to evaluate text samples on the basis of *both* reception and production.

4.2 Describing the Population

A distinction between external and internal criteria is of particular importance for constructing a corpus for linguistic analysis. The internal criteria are those which are essentially *linguistic*: for example, to classify a text as formal/informal is to classify it according to its linguistic characteristics (lexis/diction and syntax). External criteria are those which are essentially *non-linguistic*. Section 7 contains a list of attributes which we consider relevant to the description of the language population from which corpus texts are to be sampled. These attributes however are all founded upon extra-linguistic features of texts (external evidence). Of course, the internal criteria are not independent of the external ones and the interrelation between them is one of the areas of study for which a corpus is of primary value. In general, external criteria can be determined without reading the text in question, thereby ensuring that no linguistic judgements are being made. The initial selection of texts for inclusion in a corpus will inevitably be based on external evidence primarily. Once the text is captured and subject to analysis there will be a range of linguistic features of the text which will contribute to its characterisation in terms of internal evidence⁴. A corpus selected entirely on internal criteria would yield no information about the relation between language and its context of situation. A corpus selected entirely on external criteria would be liable to miss significant variation among texts since its categories are not motivated by textual (but by contextual) factors.

5 Markup

Markup, here, means introducing into the text, by means of some conventional set of readable labels or tags, indicators of such text features as, e.g., chapter, paragraph and sentence boundaries, headings and titles, various types of hyphenation, printers' marks, hesitations, utterance boundaries, etc.

5.1 Methodological Considerations

5.1.1 Converting Written Material

The Text Encoding Initiative has already published draft guidelines dealing with the markup of machine-readable texts for interchange, and proposing a markup which achieves a high level of detail and descriptive generality. The cost of introducing a sophisticated markup is high in terms of manual effort, and the cost/benefit balance may be badly upset unless it can be justified by worthwhile gains in ease of processing, accuracy and re-usability. What is needed in order to build up a large and useful corpus is a level of markup which maximises the utility value of the text without incurring unacceptable penalties in the cost and time required to capture the data.

⁴ Biber and Finegan have published widely on this subject: e.g. Biber, D. (1989) 'A typology of English texts', *Linguistics* 27.

There are two approaches which can be taken in marking-up texts. First the *descriptive*, in which one tags the underlying structural features of a text (the sentences, paragraphs, sections, footnotes, etc.). These are usually signalled by spacing, punctuation, type font and size shifts and so on, but there is no one-to-one correspondence between features and realization. Second the *presentational*, in which one tags these typographical features themselves.

The two types of markup are not mutually exclusive categories; they describe the two ends of a scale. In some instances it will be important to “record what’s there on the page”, since the researcher will be concerned to discover patterns of usage relating to features of punctuation and spelling. In other cases, the computational considerations relating to the processing of the text in indexing and free-text retrieval software are likely to take precedence and a descriptive approach will be preferred.

In the case of sources which need to be converted from printed form, markup will have to be introduced. This could be carried out during the process of OCR scanning, during keyboarding, or as a post-editing operation once the plain text has been captured. In the case of data which comes from other computer systems, there is often some explicit descriptive encoding of text structure, often arbitrarily intermingled with presentational codes, which can be parsed and converted automatically into SGML format. In both cases there is a potentially large hidden cost involved in introducing, converting and standardising markup by program, even though this appears at first to be significantly more cost effective than tedious and costly human editing. The programming work required to normalise the encoding and markup must often be repeated for each new text, because there is very little standardisation in printing format or in WP packages.

When adding a new text to a corpus under construction, one has to make the decision as to whether it is an instance of a pre-defined general text type, or whether it deserves separate treatment with different markup. That is, the SGML markup will require the specification of a set of *document type definitions* (DTDs) which formally define the structure which is to be marked up for each document type. The advantage of creating more DTDs is that the markup will better reflect the structure and content of each of a diverse collection of text samples. The disadvantage in creating more DTDs is that the marking up of a large number of different texts from many different sources becomes increasingly complex, and the generalising power of the markup is diminished. The increased complexity of the markup may not be justified, however, if most of the processing on the corpus is to be done over large aggregates of text rather than individual samples.

5.1.2 Transcription of Speech

The TEI has not yet published any guidelines relating to SGML markup for speech in transcription. One reason for this may be that gathering and transcribing authentic speech is quite a different operation from handling written documents and is clearly a much more specialist area of activity, for linguists, lexicographers and speech technologists.

This section does not address the problems of the design of acoustic corpora—corpora intended to assist in the analysis of the physical characteristics of speech—in which we have no specialist expertise. We assume here that spoken language is to be collected in large quantities in order to form the basis for quantitative studies of morphology, lexis, syntax, pragmatics, etc.

Many media organisations or research establishments will be able to supply paper or machine-readable transcriptions, which can be processed to bring them into conformance, as far as possible, with a standard markup. Transcriptions made by non-specialists will typically be in the form of quasi-written text. That is, there will be recognisable sentences and punctuation, with a high degree of normalisation of false starts, hesitation, non-verbal

signals and other speech phenomena. This type of transcription converts spoken language into a form of idealised ‘script’ (like a screenplay or drama script) which conforms to many of the established conventions of written English. Unless the corpus is intended to serve the needs of speech specialists, then the usefulness of a ‘script’ transcription is sufficient for a wide variety of linguistic studies. The advantages of transcribing in this way are:

- the cost and time of transcription are minimised
- the transcription is easily readable without any special training
- the transcription can be processed using established and widely available text processing software without substantial pre-editing

5.2 Features for Markup: written text

Non-ascii characters: The number of such characters which may need encoding could be quite large. These might be encoded at the time of data capture and for the sake of economy, might be recorded using any local convention which ensures that the codes are unambiguous yet easily keyed and checked. They can be expanded into standard SGML entity references by a search-and-replace operation carried out later. In many cases texts which are received already in machine-readable form will include special codes for graphic shapes which are not specified as part of the ascii set. These will have to be identified and standardised.

Quotation: This is an important aspect of text encoding in a corpus. There are three types of quotation that need to be considered:

- direct speech
- block quotes
- other uses of quotation marks

Direct speech is probably the most difficult of these to deal with satisfactorily. In order to tag or parse direct speech accurately, it would be necessary to distinguish the multiple levels of *subsidiary* and *primary* discourse, each with its own syntactic structure. Block quotes are much more easily handled. Either in print or in machine readable form, most texts signal the start and end of block quotations (with indentation, typeface change, opening and closing quote marks, etc.) Such representational features need not be recorded, as long as the extent of the quotation is indicated. Other uses of quotation marks, to signal ironic or jocular uses, cited words, titles of books and films, etc., may be marked with SGML entity references for the opening and closing punctuation marks.

Lists: If lists are not marked up in any special way they give rise to a number of undesirable side-effects in text analysis. First, the *item labels* of a list may be roman or arabic numerals, letters or other printer’s mark, and these can be misinterpreted by text searching software as ‘real’ words. Second, lists are often not punctuated according to the normal conventions with respect to sentences. This is likely to confuse tagging and parsing software.

Headings: Headings can usually be interpreted as labels attaching to some structural unit: chapter, section, article, and so on. Or they could be treated as short interruptions in the main flow of the text. Unfortunately, real-life texts are much less tidy than our idealisations of them, and often texts are found in which apparent headings do not seem to be functioning as labels to any structural unit. Newspapers and magazines are particularly inconsistent in this respect.

Abbreviations, initials and acronyms: A surprisingly high proportion of the word tokens of a large English-language corpus will be accounted for by abbreviations, initials and acronyms, e.g. personal names, organisations, titles of address, postcodes, units of measurement, days of the week, month names, chemical elements, conventional Latin-derived abbreviations. Of these a substantial proportion could be identified automatically by pattern matching and tagged as contracted forms. Automatic identification and tagging of abbreviations and initials in a corpus may be supplemented by manual editing with the aim of eliminating most of the overlap between words and abbreviations.

Front and back matter: Books will usually include a certain amount of front matter (e.g. preface, foreword, contents, list of figures, acknowledgements) and back matter (e.g. index, appendices, bibliography). Some of these may be captured and included in the corpus, while others may be omitted.

Chapters and sections: These can easily be encoded with little extra effort and keying. Manual editing of the encoding of major text divisions will not be too time-consuming.

Proper names: The aim, in marking these, is to resolve by pre-editing the ambiguity between proper names and other homographic word forms. If carried out manually this is a major undertaking. Word class tagging software can be used to identify a large percentage of proper names automatically.

Correspondence and addresses: The conventional paraphernalia which attach to the body of a letter (addresses, dates, ‘cc’ lists, salutation, document references, etc.) need to be handled in an appropriate way to ensure that the non-discursive material is identifiable automatically by processing software.

Pagination: Generally, it will not be possible in all cases in a large corpus to preserve the pagination and numbering, especially if texts are acquired in electronic form. It may however be important to users of the corpus to be able to refer to the actual page of the printed original of the text.

5.3 Features for Markup: spoken text

Speaker change: The basic structure of speech transcription should be a sequence of speaker *turns*. Each turn should begin with an encoding identifying the speaker wherever possible. The use of a minimal encoding will reduce the amount of keyboarding required at the data capture stage.

Syntax: It will be a matter for each corpus project to decide whether, and to what extent, punctuation and markup reflects a notion of normal syntax. Word-class taggers and parsers are currently oriented almost exclusively towards analysing written language⁵ and it may be important to preserve as far as possible the syntactic units of clause and sentence. Other projects may prioritise prosodic analysis of the text and use a markup which encodes tone units or other segmentation strategies.

Accent, dialect and standard orthography: It is important for the corpus builder to decide to what extent the transcription is to represent the sounds of speech (e.g. accent) and to adopt a transcription encoding which is adequate for the representation. If a plain prose-style transcription is adopted, the transcription should be consistent

⁵ Svartvik and Eeg-Olofson have developed a tagger and parser tailored to analysing speech in the London-Lund corpus.

in the use of orthographic irregularities and clear specifications should be given for the use of enclitic forms, variant spellings, and the use of non-standard spelling forms which might be used by the transcriber to reflect the sounds of speech. (E.g. *don't*, *can't*, *yeah*, *dunno*, etc.) A closed set of permissible forms can be given to the transcribers for guidance. Similarly, one must consider whether dialect forms are to be preserved in the transcription or regularised to conform to an agreed standard. Since it is probable that a large natural language corpus will be valuable for the study of dialect forms (in syntax and lexis) the standardisation should probably not be applied so strictly that these distinctions are lost.

Interruption and overlapping speech: The writing system does not have very well-established conventions for representing this feature of speech. Interruption is a feature which can without difficulty be represented in the linear stream of writing. It merely requires the insertion of a code or tag which indicates that the preceding turn is incomplete because of the intrusion of the following turn. Interruptions are sometimes more messy, however, and the interrupted speaker may choose to continue regardless of the rival turn: the result will be overlapping speech, which cannot be so easily encoded in linear written form. A detailed markup would encode for each overlapping segment of speech

- the start and end points
- an identifier
- the speaker
- the point at which this segment begins overlapping with another
- the point at which this segment ends overlapping with another
- the continuity of each speaker's turn

The markup could become very dense and specialist training and skill could be required to carry out transcription if precise details of overlapping speech are to be recorded faithfully. A sophisticated markup would require substantially more time and effort than a simplified encoding, and the cost of staff training, quality control and the additional time required for analysis and keying should be carefully estimated before final decisions are made.

Pauses: Pauses may be voiced or silent, long or short. It is very simple for a transcriber to record the voiced pauses and they can be encoded in several ways. SGML entity references might serve this purpose quite well. Silent pauses are more problematic, because a transcriber who has only an audio recording of the speech event cannot be sure what other activities or interference might be the cause of a silence on the tape. Unless the recordings are analysed in detail in relation to the physical action of the speech event, the encoding of silent pauses is likely to be misleading and unhelpful. Voiced pauses can be encoded using two codes; one for a short noise and another for a long one.

Functional and non-functional sounds: Some speech sounds have a clear discourse function. The recognised functions may be simplified to a small set each of which could be represented as a standard code. This places responsibility on the transcriber to interpret speech sounds and assign them to appropriate functions. Alternatively one could devise a large set of orthographic representations to cover a full range of speech sounds and encode the sound rather than its discourse function. If the corpus is primarily for grammatical and lexical studies, then the precise recording of functional speech sounds will not greatly enhance its value. Laughter, coughs, grunts, and other sounds which may occur in recording can be simply encoded.

Inaudible segments: Often, especially if recordings are made in real-life situations, extraneous noise or poor recording conditions will make it impossible for the transcriber to hear exactly what is said. Lacunae should be marked with an indication of the extent of the inaudible segment (measured roughly in, say, seconds, syllables, or ‘beats’ of speech rhythm)

Spelling: Some words may not be familiar to the transcriber and cannot be spelled with certainty. Proper names and technical terminology are likely to be especially difficult for a transcriber. An attempted spelling can be made at the time of transcription and a marker inserted to allow correction during post-editing, or at least to indicate to the corpus user that the spelling is doubtful.

Numbers, Abbreviations, Symbols: The use of alternative written forms for such words as *Mister (Mr.)*, *pounds (£)*, and *two hundred and thirteen (213)* needs to be considered. Since there is no sense in which the original form can be faithfully reproduced, a policy of normalisation might be followed to make subsequent processing easier.

6 Corpus Typology

A corpus is a body of text assembled according to explicit design criteria (see Section 7 below) for a specific purpose, and therefore the rich variety of corpora reflects the diversity of their designers’ objectives.

It is worth mentioning in parenthesis that the text typology discussed in Section 7 is, in many cases, also relevant to corpus typology, in that corpora may be classified according to text types if they consist solely of texts of one single type. Thus, if the corpus is created for the purpose of studying one single MODE, then one may have a SPOKEN or a WRITTEN corpus; similarly, if only one MEDIUM is of interest, one may have A BOOK or a NEWSPAPER or a CLASSROOM LESSON corpus.

In this section, however, our purpose is to outline certain contrastive parameters of corpus typology per se:

1. **Types:** FULL-TEXT SAMPLE MONITOR

Notes: For Full Text: each text in the corpus is unabridged.

For Sample: sample size to be defined, also location of sample within full text and method of selection of samples.

For Monitor: texts scanned on continuing basis, ‘filtered’ to extract data for database, but not permanently archived. (Clear(1988), Sinclair(1982))

2. **Types:** SYNCHRONIC DIACHRONIC

Notes: A specific period must be designated for a synchronic corpus; this requires research into how long that period may be if the corpus is to be considered synchronic.

3. **Types:** GENERAL TERMINOLOGICAL

Notes: Terminologists must define conditions which must obtain if a corpus is to be valid for terminological use, this in terms no doubt of the text typology (see Section 7 below)

4. **Types:** MONOLINGUAL BILINGUAL PLURILINGUAL

5. **Types:** LANGUAGE(S) OF CORPUS

Notes: English, Russian, Japanese, ...

6. **Types:** SINGLE PARALLEL-2 PARALLEL-3 ...

Notes: Are all the texts in the corpus stand-alone or part of a parallel pair/trio etc. of translated texts? (This applies only to bi- or plurilingual corpora)

7. **Types:** CENTRAL SHELL

Notes: The central corpus is a selected body of texts, of manageable size, big enough for normal purposes. The shell, which may be the remainder of the ETL, is available for access when necessary.

8. **Types:** CORE PERIPHERY

Notes: These concepts are discussed by Leitner⁶ in relation to the ICE. The ‘core’ contains text types common to all varieties of English, and therefore present in all the subcorpora; while the ‘periphery’ contains text types specific to some subcorpora only.

7 Text Typology

There is much talk of a ‘balanced corpus’ as a sine qua non of corpus analysis work: by ‘balanced corpus’ is meant (apparently) a corpus so finely tuned that it offers a manageably small scale model of the linguistic material which the corpus builders wish to study. At present corpus ‘balance’ relies heavily on intuition, although work on text typology is highly relevant. It is not our purpose to lay down a methodology for ‘balancing’ a corpus. For we believe that this is not something that can be done in advance of the corpus-building process, if it can be done at all. Controlling the ‘balance’ of a corpus is something which may be undertaken only after the corpus (or at least an initial provisional corpus) has been built; it depends on feedback from the corpus users, who as they study the data will come to appreciate the strengths of the corpus and be aware of its specific weaknesses.

In our ten years’ experience of analysing corpus material for lexicographical purposes, we have found any corpus—however ‘unbalanced’—to be a source of information and indeed inspiration. Knowing that your corpus is unbalanced is what counts. It would be short-sighted indeed to wait until one can scientifically balance a corpus before starting to use one, and hasty to dismiss the results of corpus analysis as ‘unreliable’ or ‘irrelevant’ simply because the corpus used cannot be proved to be ‘balanced’. It should also be noted that recording attributes of texts is very labour-intensive and an over-ambitious system could strangle a corpus at birth. Experience teaches that it is better to aim to record initially an essential set of attributes and values which may later be expanded if resources permit.

The significant variables considered here, in the context of corpus design criteria, are all extra-linguistic. We believe however that it is impossible to ‘balance’ a corpus on the basis of extra-linguistic features alone. Diagnosis of imbalance must come from an analysis of internal evidence. All that the corpus-builder can do is to try not to skew a corpus too much in any direction. Balancing it, or at least reducing the skew factor, is something which comes along much later, and will demand information on both linguistic and extra-linguistic features in the corpus.

When creating a corpus for a specific purpose, the corpus designer must be able to make principled choices at all the major decision points. Information on features fundamental to corpus design must therefore be recorded on each of the texts in the electronic text library from which the corpus is to be selected. We suggest that the concept of a set of features

⁶ in Section 4.2 of ‘Corpus design—problems and suggested solutions’, ICE working paper, May 1990

is relevant to an ETL, while that of a taxonomy proper relates to a corpus created for a specific purpose. See Francis & Kucera (1964), Johansson, Leech and Goodluck (1978), Leitner (1990), Oostdijk (1988), Engwall (forthcoming).

Looking at features relevant to the typology of texts, we identify extra-linguistic variables which are of interest to anyone establishing an ETL or a corpus; we propose that these variables may be considered criterial attributes in the context of ETL/corpus design, and we try to indicate the minimum level of detail which seems to us essential to record. We exemplify the set (sometimes open-ended) of values of these attributes, and for some consider the type of criteria that may be applied in the process of identifying the features of a text in an ETL or a corpus.

7.1 Text Attributes

This section consists of:

1. a list of attributes that may be recorded for every text introduced into the text collection
2. a note regarding criteria for assignment of values;
3. an example of acceptable values and an indication of a reasonable level of depth to be recorded; values in capitals are, we believe, essential to note; those in Roman type are highly desirable but not essential in the first instance, where insistence on this degree of detail could be counter-productive;
4. where appropriate, some notes on the value-assignment criteria.

1. **Attribute:** MODE

Criteria: mode of delivery of the original contents of the text: text transcribed from an audio/video recording is ‘spoken’; text written as dialogue is ‘written to be spoken’; poetry poses a problem (written to be read or spoken?)

Values: WRITTEN (default)
written-to-be-read (default)
written-to-be-spoken
SPOKEN
spoken-to-be-written

2. **Attribute:** PARTICIPATION

Criteria: self-evident: number of people originating the text

Values: 1-PERSON (default)
2-person
multi-person

3. **Attribute:** PREPAREDNESS

Criteria: Still to be defined.

Values: PREPARED (default) (written; never spoken)
scripted (spoken, never written)
from-notes (ditto)
spontaneous (ditto)

Note: Some of the points on this cline are linked to the attribute of MODE in such a way as to make the PREPAREDNESS attribute often redundant, ie : if WRITTEN, then always PREPARED; if SPONTANEOUS, then always SPOKEN.

This parameter subsumes the notion of timed vs non-timed.

4. **Attribute:** MEDIUM

Criteria: medium of original contents of the text; all written text, whether published or private, is TEXT; the others are mostly self-evident, though criteria for assigning values to problematic cases such as films (scripts or soundtracks) need clarification.

Values: (*Writing*)

TEXT (default)

BOOK

PERIODICAL

NEWSPAPER

typescript

manuscript

(*Speech*)

distant direct

TV talk (lecture/speech, etc)

radio entertainment (theatre, etc)

telephone person-to-person

Note: Linked in one place to the attribute of MODE (ie if TEXT, then always WRITTEN).

5. **Attribute:** STYLE

Criteria: mainly self-evident

Values: PROSE (default)

poetry

blank verse

rhyme

6. **Attribute:** GENRE

Criteria: mainly self-evident

Values: (*Writing*)

NOVEL

SHORT STORY

PLAY

POEM

ESSAY

LETTER

business

personal

(*Written or Spoken*)

advertisement

regulation/law

article

advice column

horoscope

announcement

(*Speech*)

lecture

debate/discussion

speech

conversation

demonstration

classroom lesson

examinations

report

commentary

feature

advice programme

Notes: Some problem areas (how long is a short story? how much of TV news programmes is report and how much commentary? etc). Another big problem with this attribute is attaching values to the very small pieces of text that go to make up the whole of a newspaper or periodical (very labour-intensive); but see 1.2.1 above.

7. **Attribute:** CONSTITUTION

Criterion: Self-evident: is the text single or composite? A single text by one author is single; a newspaper, journal, collection of essays, textbook etc is composite, ie made up of many distinct small texts (which could each be classified individually).

Values: SINGLE (default) COMPOSITE

8. **Attribute:** FACTUALITY

Criteria: To be defined: many problem areas—how to mark history, biography, autobiography etc?

Values: FACT faction FICTION (cline)

Notes: Closely linked to attribute of GENRE, eg the value NOVEL or SHORT STORY there implies FICTION here, etc.

9. **Attribute:** SETTING

Criteria: In what social context does the text belong?

Values: social educational
 business

10. **Attribute:** FUNCTION

Criteria: To be defined. Not easy.

Values: UNMARKED (default)
 narrative
 informative
 hortatory/persuasive
 regulatory/instructional
 reflective (=giving one's opinion) etc
 creative/artistic

11. **Attribute:** TOPIC

Criterion: "This text is about X", or "The subject of this text is (related to) X"

Values: science music
 biology orchestral
 chemistry opera
 etc etc

Note: It is necessary to draw up a list of major topics and subtopics in the literature.

12. **Attribute:** TECHNICALITY

Criteria: Based on degree of specialist/ technical knowledge of the author and target readership / audience

Values: GENERAL (default)
 (non-specialist author & target)
 TECHNICAL
 (specialist author & specialist target)

semi-technical
(specialist author, general target)

Note: These must be external variables, not linguistic style variables. This particular attribute is highly important for terminological corpora.

13. **Attribute:** DATE

Value: DATE OF PUBLICATION (default) or date of speech event

Note: this is related to 14. TEXT STATUS. For some corpora it will be desirable to note the date of first publication, in case of revised editions.

14. **Attribute:** TEXT STATUS

Criterion: Is the text in its first appearance? Or a reprint? Or a 'new' or 'revised' or 'updated' edition?

Values: ORIGINAL/REPRINT (default)
updated
revised etc.

Notes: This attribute is of particular importance in terminological corpora but also of interest in literary corpora.

15. **Attribute:** LANGUAGE

Criterion: Self-evident

Value: English, French, Japanese ...

16. **Attribute:** LANGUAGE LINKS

Criterion: Self-evident: is the text stand-alone or part of a parallel pair/trio etc of translated texts?

Values: SINGLE (default)
PARALLEL-2
PARALLEL-3 (etc.)

17. **Attribute:** LANGUAGE STATUS

Criterion: Is the text the source language or is it a translation?

Values: SOURCE (default) TRANSLATION

Note: Relevant only when the value of 16. LANGUAGE LINKS is not single.

18. **Attribute:** METHODOLOGY FOLLOWED

Criterion: To be defined: depends on assessment or knowledge of whether this is concept-based, whether it is standardized or officially approved, what research principles were applied.

Values: To be determined: possibly on a scale from 1 to 5?

Notes: This attribute of importance in terminological corpora. The criteria and values must be discussed with a terminologist.

19. **Attribute:** AUTHORSHIP

Value: NAME OF AUTHOR(S)

20. **Attribute:** SEX OF AUTHOR(S)

Criterion: self-evident

Values: MALE FEMALE

21. **Attribute:** AGE OF AUTHOR(S)

Criterion: self-evident

Values: the actual age in years

22. **Attribute:** REGION OF AUTHOR(S)

Criterion: To be defined: this attribute relates to the regional type of the language of the author(s)

Values: STANDARD (default) (for those languages where applicable)

(other values will depend on specific language: eg for English—UK, USA, Canada . . . , or English, Scottish, Welsh . . .)

Notes: Is a brand of language to be described as (e.g.) ‘Irish’ if the AUTHOR(S) spent the first 20 years of life in Ireland and the last 50 in the US? etc etc This is a parameter which will be refined by internal evidence.

23. **Attribute:** NATIONALITY OF AUTHOR(S)

Criterion: Mainly self-evident.

Values: actual nationality at time of writing/speaking

24. **Attribute:** AUTHOR(S)’S MOTHER TONGUE

Criterion: Self-evident

Values: actual language, if known
default = language of text

Notes: This attribute of particular importance in terminological corpora.

25. **Attribute:** AUTHORITY OF AUTHOR(S)

Criteria: To be defined. This relates to credibility of author(s) and authority on subject-matter; determination of this will depend on qualifications, experience etc.

Values: Probably point on scale from 1–5.

Notes: This attribute of importance in terminological corpora.

7.2 Reference Coding

The set of data which is to be recorded for each text sample included in a corpus may well become fairly large. In addition to the attribute values presented above, it may include discursive information about the editorial decisions taken when the text was captured and validated or other information which will help future users of the corpus understand the precise nature of each sample. This data can be conveniently stored separately from the actual textual material of the corpus, but this will require an adequate reference system to relate any word in the corpus back to its location in the original text and to the associated corpus information. If the corpus is made up entirely of published printed texts, then a reference system based upon standard bibliographic citations and page numbering will probably be appropriate. However, for transcribed speech, printed ephemera, manuscript material, private or business correspondence and similar material some other method must be used to refer from the electronic form of the corpus to the original. In some projects this may not be necessary and a simple text referencing system may be adequate. In others where, for example, audio and/or video recordings accompany the computer corpus texts, a

sophisticated encoding will be necessary to allow the user to consult the original or determine exactly where a given stretch of discourse is located. For many corpus applications sentence number will be a suitable identifying unit within each text sample.

7.3 Case Study: Designing a Lexicographic Corpus

The previous section sets out a range of text classification features which would allow users of the corpus or ETL to extract subsets of the information in many different configurations according to the focus of interest. This section describes an approach to the classification of texts for a particular corpus at Oxford University Press which is being constructed primarily to meet the needs of lexicographers and grammarians working on both native-speaker and ELT reference works.

First a simple taxonomy of text types is described which is intended to cover the broad range of modern English language both spoken and written. This taxonomy takes no account of the *topic* or *subject area* of the text, which is to be treated as a second dimension of text categorisation. The third classification is based on a set of features (some binary, others graded) which are independent of the traditional categories presented in the taxonomy.

7.3.1 A Taxonomy of Text Types

Spoken

Dialogue

Private

Face-to-face conversation

Structured

Unstructured

Distanced conversation

Classroom interaction

Public

Broadcast discussion/debate

Legal proceedings

Monologue

Commentary

Unscripted speeches

Demonstrations

Written

To be spoken

Lectures

Broadcasts (news, documentary)

Drama scripts

Published

Periodicals

Magazines

Newspapers

Journals

Newspaper supplements

Books

Fiction
Non-fiction
 General
 Official reports
 (Auto-)biography
 Reference (discursive)
 Educational textbooks

Miscellaneous
 Brochures
 Leaflets
 Manuals
 Adverts

Unpublished
 Letters
 Personal
 Business
 Memos
 Reports
 Minutes of meetings
 Essays
 Other

Notes

- This classification is *not* exhaustive.
- The classification of spoken text is taken, with only minor modifications, from Greenbaum’s specifications for the International Corpus of English.
- The ‘overlap’ category of text which is written to be spoken is included within the class of written material. It could be treated within the spoken class, since the actual performance is likely to manifest characteristic features of speech (hesitation, false starts, syntactic loose ends, etc.) or indeed it could be treated as a special primary class in its own right.
- Structured face-to-face conversations include meetings and seminars in which the turn-taking is generally free but the interaction is guided explicitly (by a chairperson) or by some sort of acknowledged agenda. Unstructured conversations are not constrained in this way.
- The books/fiction class is not further subdivided at this stage but may be later.

7.3.2 Topic

The subject field (hereafter, ‘topic’) of a text is difficult to deal with in any classification. First, it is not clear to what extent the topic is an external or internal factor. Newspaper reportage, for example, seems to be *prima facie* on the topic of current affairs; it is not necessary to read the text in question to be able to classify it thus. On the other hand, the ‘aboutness’ of a text is substantially a function of its lexis and is therefore an internal matter—a text is ‘about’ nuclear physics if it uses words and phrases which refer semantically to the discipline of nuclear physics.

Second, labelling the topic is a classification of the world, and the world is a very complex structure which does not submit easily to straightforward classification. We have

developed a conventional set of terms to designate aspects of the world of learning and activity: science, arts, social science, humanities, technology, crafts, etc. Library science has developed classification systems to meet the basic needs for information retrieval. A complete list cannot be enumerated, since the world can be categorized in infinite detail. Moreover, it seems likely that beyond the crudest labelling, the organisation of topics resembles a network more than a hierarchy. Nevertheless, a system based on library cataloguing may prove to be most appropriate for the purposes of this particular corpus.

Third, these familiar labels cut across the text types identified in the preceding section. A business letter, for example, may be classified with any topic label. Although this is gross simplification it does seem that there is some basic form versus content distinction here: text types define the form, while subject field defines the content.

As a strategy for corpus compilation, for each text that is held in machine-readable form as part of the corpus collection we can record its topic for subsequent retrieval and analysis (if we wish to study a subset composed of texts relating to Business and Finance, perhaps) and review the spread of topic labels periodically. In order to avoid any gross imbalance in the appearance of specialized text in the corpus, we should define some macro categories (current affairs, science and engineering, social science, leisure-sport, leisure-arts, etc.) and ensure that the sampling across these is controlled.

The topic classification can be laid on top of the text types presented above, with some exceptions. Newspapers, Newspaper supplements, Fiction, Biography and Personal Letters are not susceptible to topic classification in the usual way. These are difficult classes and many have no clear subject field. Newspapers and supplements are composite documents (made up of many discrete texts) which range over a large number of topics.

7.3.3 Controlled Parameters for Written Texts

In this section we propose a minimal set of parameters against which written text can be evaluated. These parameters are primarily *external* and should be applied as such. It is intended that each of these parameters is, as far as possible, independent of the others. These parameters will, in some cases but not all, reflect the text type classification, and may therefore appear redundant. However they will serve to ensure that a range of variation is represented independently of text type categories.

It is not intended that texts will be selected in order to fill all the possible configurations in equal quantity. Certain configurations of these controlled parameters will therefore tend to predominate. The value of classifying in these terms is that such predominance or under-representation can be clearly identified and corrected if this is felt to be necessary.

Intended Readership

Size A scale.

- 0 one person
- 1 small number of people (less than approx 20)
- 2 approx 20–100 people
- 3 large but homogeneous (e.g. within a company, college, etc.)
- 4 large geographically defined (e.g. town or region)
- 5 unrestricted heterogeneous

Education A scale from 0 (basic literacy) to 5 (university doctoral level).

Known/Unknown A scale.

- 0 very familiar (close friends, family)

- 1 familiar (colleagues)
- 2 known by name (acquaintances)
- 3 group defined by lifestyle, profession, interests, etc.
- 4 unknown except by approximate age, region, etc.
- 5 completely unknown

Age A binary choice between adult and non-adult. Adults are, say, over the age of 18.

Authorship

Sex Discussion continues on the value and practicality of recording the sex of the author(s).

Nationality and Domicile Labels such as ‘British English’ and ‘US English’ need to be defined in a practical way to take account of authors’ situation.

Number A text is not always the product of one individual. We can classify the text by the following labels:

- single
- multiple
- corporate
- unclassified

Format

Length Expressed as a word count.

Constitution *Single, composite or collection*. E.g. newspapers are composite; a number of magazine advertisements treated as one text is a collection.

Explicit structure A scale from 0 (unstructured) to 5 (highly structured). A novel is a long text with typically very little explicit structure—the chapter is often the only unit of organization. Other texts are structured into sections, articles, abstracts, summaries, numbered paragraphs and so on. This explicit structure may well need to be encoded using SGML markup.

8 Progress to date with Standards for Corpora

8.1 Terminology

Standards in respect of corpora mainly concern the compatibility of the kinds of annotation used for texts. These we shall call ‘encoding standards’. However, there is also in principle an issue as to the comparability of different corpora, including perhaps judgements on their suitability for different tasks. These we shall call ‘evaluation standards’.

8.2 Encoding Standards

There is now a fair, and growing, number of corpora becoming available which attempt to provide basic grammatical tagging of the texts they contain. It is not difficult to find differences in the regimes adopted, and hence to recognize the scope for arbitrary incompatibility in what is not a very controversial area (at least by comparison with syntactic parsing, or semantic tagging and analysis).

Any particular comparison speedily becomes lost in a thicket of details. However, it is possible to distinguish a number of different axes on which disagreements are possible, and perhaps this is the first step towards resolving them.

We may ask then:

- which levels of constituent are taggable (from character —e.g. punctuation marks— and morpheme through word to phrase, sentence and higher units);
- whether the tags are atomic in form, or whether they are effectively complexes of features;
- whether tags are assigned to some singleton classes of words, or these are left to represent themselves (e.g. *the*);
- whether there are formal constraints on assignment of tags (especially, whether non-contiguous items can share a tag);
- whether in some cases the tagging system allows free variation in tag to be assigned— (e.g. is ‘male frog’ N+N or A+N or both?).

All these questions have been answered differently by the compilers of corpora currently available; no doubt there is scope for more variation too. In many cases decisions have been provoked by particular applications and research goals; in others by physical constraints of the amount of data to be processed within a given time, or the technical facilities available.

As corpora are more and more seen as a fundamental research tool, however, for multiple uses, it becomes worthwhile to attempt explicitly to provide a common framework. The objective is, initially, as far as possible to define existing annotation schemes in terms of this framework; ultimately, to add substantive details to it so as to create an all-purpose notation scheme, with resources adequate for all the major applications. Since new applications will no doubt arise in future beyond those currently foreseen, it will also have to have clear rules for extension.

The Text Encoding Initiative (TEI) of the Associations for Computers and the Humanities (ACH), for Computational Linguistics (ACL), and Literary and Linguistic Computing (ALLC) is the only attempt known to the authors deliberately to propose a common standard in this field. It is notable, however, that its Guidelines as so far issued (July 16, 1990), (which are an extension of SGML, the Standard Generalized Mark-up Language for electronic text formatting,) do not go beyond the high-level goal of providing a syntax of labelled bracketings. Hence an editor attempting to provide a set of texts in accord with them is still compelled to make a large number of decisions of his own. (Cf, e.g., Liberman’s commentary on his redaction of texts for the ACL Data Collection Initiative.)

It is asking a great deal to expect a standard notation to provide a usable common framework for the actual categories used in grammatical tagging, even supposing agreement is reachable on the more formal considerations mentioned above. However, at this point it becomes possible, in principle, to benefit from other related standardization work, viz that undertaken to identify and promote reusable lexical resources, and to set up compatible formats for machine-readable dictionaries.

The grammatical tags used in corpus annotation will need largely to convey the same information that is assigned to lemmata in dictionaries. Evidently, where corpora are to be processed by NLP devices that access machine-readable dictionaries, it will be simplest if the tags are identical with the dictionaries’ grammatical categories.

There are two projects in Europe which are addressing this problem.

ESPRIT ACQUILEX (to finish in mid 1991) addresses itself to machine-readable dictionaries as currently available (e.g. as typesetting tapes) for a number of European languages (viz English, Dutch, Italian), monolingual and bilingual, and aims to define a useful common notation for their grammatical categories.

EUROTRA-7 is a feasibility study, rather than a concrete project, due to end in April 1991: it will survey existing terminological and lexical resources of all kinds, and assess the feasibility of standardizing them. It will include both monolingual and bilingual lexical resources, and consider possible architectures for a system that would make them available for potential re-use.

In general, one can expect advances in developing the tag-set on a standard basis to occur step by step with standardization of dictionary coding; and it is largely to the lexicographer that corpus-builders will have to look, for expertise in devising consistent but practical schemes.

It is clear that we are just taking the first concrete steps towards setting up the standards that will make the linguistic annotation of corpora compatible with arbitrary processors. At least the requirements, and to some extent the capabilities, of linguistic processors are clear, since there is a coherent natural language processing community world-wide, now supplemented with forward-looking lexicographers.

However, once the wider uses of corpora identified in Section 9 begin to become real, the need for more general standards will also become evident. These may require annotation codes for a wide variety of subject-matter, and will also involve many differing communities of users. It is to be hoped that, by that time, progress made in defining workable standards for linguistic processing of corpora will have provided insight into how best to extend coverage to quite different fields.

8.3 Evaluation Standards

Evaluation standards are less of an issue in respect of text corpora than they are for many other constructs in natural language processing. There is in fact little danger of obfuscation for the major parameters that characterize a corpus: its size (in numbers of running words), and gross characterizations of its content.

It is necessary to keep sight of both of these parameters when comparing and evaluating results derived from different text corpora. Mutual information figures are only comparable as between items derived from the same sizes of corpus. And evidently it is always at least possible that any statistical relationship observed is crucially dependent (however indirectly) on the choice of subject-matter that occurs in the corpus.

Having said this, however, it becomes clear that the major issues in this sort of standardization are still waiting for the resolution of current disputes. In order to get an uncontroversial view of what type of corpus evidence is required substantively to answer questions of a given type, we must first find agreement on whether ‘balancing’ of a corpus is ever possible, let alone ever necessary. And that itself is likely to have to await the accumulation of empirical results.

It would seem, then, that evaluation standards for corpora are still premature.

9 Potential Users and Uses

9.1 Generalities

The large files that constitute a corpus, an electronic text library or a text archive have come together as if by chance. They are together because of the language they are written in, and it is very unlikely that this had anything to do with the author’s motive in writing a text.

This means that users of these files are unlikely to be interested in the content of any particular text. Texts appear as specimens, and users will be interested in what they show

about the class represented, not in any of the particular information the text was created to impart.

This makes corpus-users rather untypical among the vast ranks of people who are interested in consulting computer files. In fact, corpus-users can be divided into three types: those interested in the language of the texts, those interested in the content of texts (as representative of a larger collective), and those interested in the texts themselves as a convenient body of test material for electronic media.

9.2 Language Specialists

This class will include the earliest users of corpora. At the moment the most active groups supporting corpus development are drawn from this class, though the picture is slowly changing, especially in Japan.

Lexicographers They will use the corpus for information on the actual usage of the words they cover. It may be consulted directly, for information on specific words; or it may be processed in various ways, in order to develop parts of a lexical database. The processing is likely to require any or all of the electronic processing media surveyed in 9.4 below.

Language Workers Translators, terminologists and technical writers will be concerned with corpora for special purposes. Parallel corpora (of texts and their translations) are beginning to emerge, especially in the EC where large quantities of documentation is prepared in several languages. Specialized corpora representing technical areas can provide the basis for the enhancement of terminology databases and contribute to the success of efforts towards standardization of technical terminology.

Computational Linguists At present these researchers separate into two camps, the 'self organizing' and the 'knowledge based'. The former attempts to use the statistical regularities to be found in mass text as a key to analysing and processing it; while the latter brings in the results of linguistic theory and logic as the foundation for its models of language.

The self-organizers use corpora first for initial tests for the presence or absence of regularities they expect to discover. Then, having created search and processing programs, they use corpora to 'train' them, i.e. to refine them through a repeated process of trial and error.

Adherents of the knowledge-based approach have only recently come to recognize the usefulness of corpora in their work, primarily to assess the level of coverage achieved by processes that they have designed a priori.

A major promise of corpus-based studies is to suggest how to integrate the work of these two groups. In particular, recent studies have suggested that self-organizing statistical techniques gain much in effectiveness when they act on the output of grammatically-based analysis. If this promise is borne out, then we can expect the self-organizing techniques to contribute much more to the semantic, rather than the grammatical or syntactic, analysis of texts: i.e. as a means of analyzing their content, rather than their linguistic form.

At the moment, such studies largely restrict themselves to showing how statistical analysis (using mutual information etc.) can give results which seem intuitively justified in any case. The next stage is to use the techniques so as to reveal semantic regularities hitherto unsuspected; and also to begin to build up general mechanisms

for the automatic content analysis of texts. The former of these will use corpus material as its main data-source; the latter will need copious corpora for testing purposes during development—though of course ultimately it will be applied to particular texts for which a real analysis is required.

Theoretical Linguists They view corpora as a mass exhibit in extenso of the facts of a language, yielding data on the relative frequency of phenomena of all kinds, and in particular providing a check on the evidence of their own, or their informants', intuitions.

Applied Linguists In teaching a second language, corpora provide a substantial resource for extracting and studying authentic language data with the authority of attested use. This data might be presented directly to students for classwork or may underpin the preparation of teaching materials of every kind. Increasingly, computer corpora and a range of software packages and tools are available as part of the language teacher's resource pool. Subcorpora of restricted topic areas may be particularly appropriate for use in teaching Languages for Specific Purposes.

9.3 Content Specialists

Corpora, as they grow so as to include large subsections classified by date, subject-matter, region, age-group or whatever, will become interesting as a data-source on the classes of people who created the texts as well on the texts' language. Such people will include historians, literary critics and sociologists, perhaps also advertisers and pollsters. At the moment, these groups are limited in their analysis of corpora to looking at the incidence of fixed words or phrases. But as the language specialists make progress, we can expect corpora to become increasingly useful to those who want a means of handling mass text files, so as to draw out trends and overall analyses.

Historians will be able to track the development of opinions and ideas through study of the words and phrases that refer to them. An example might be a historical study of the use of noun-phrases referring to machines as the subject of action verb-phrases, as indirect evidence on how our mechanistic conception of people and animals has established itself. They will also be able to use dated and analyzed corpora to discover implicit time-stamps and place-stamps, which they can then apply to identify documents whose origin is unknown.

Literary critics have already made signal use of corpus-based research under the heading of stylometrics. Statistical analysis of word-use is already crucial in determining ascription of dubious work to a known author, and such techniques can only become more effective as linguists discover the significant features which are present at a higher level than individual words.

Besides acting as a mass training ground for these techniques, corpora will also be useful as sources for statistical information on the differences of style characterizing different groups, whether by age-group, period, country of origin, or whatever. Whole areas of literary analysis which are now inescapably subjective will begin to allow some objective test—though this is unlikely to make them any less controversial!

Sociologists will be able to make similar uses of the corpus resources, though here the parameters of interest will be different: not period, author or genre, but class, race, creed.: whatever, indeed, the corpus architects may have chosen as significant labels for texts, and any combination of these labels. We can expect statistical confirmation of the whole host of perceived nuances by which people's language betrays their origin, but also the discovery of a host of others, hitherto too subtle or abstract to be picked up.

9.4 Media Specialists

Corpora will be the indispensable testbed for all the text-processing functions that software developers devise in the coming years. Inevitably, the success of these new developments can only follow on substantial progress in the research of the language specialists we have already considered. However, even now, it has been recognized by the US Defence Advanced Research Projects Agency that a large and common source of textual data is a *sine qua non* of progress both by the researchers and the developers who follow in their wake.

These text processing systems will be quite various, and there may come a time when human authors are no longer the only source of readable text, nor human readers the only users of it. At the moment, attention is concentrated above all on three types of application: Information Retrieval Systems, Machine Translation, and Speech Processing.

Information retrieval systems are themselves extremely varied. They include mechanisms to extract information that fits a given format from bodies of text, (which may be fixed or may be dynamically accumulating as real-time messages), then using it to build up a knowledge-base; mechanisms to find enough information in items (say, messages or articles) to decide on an appropriate addressee (i.e. message-routing, document-clipping); mechanisms to find the information for an index, or more ambitiously, to summarize the important content.

Machine translation (or machine-aided translation) is a computer application whose attraction, where it is feasible, is self-evident. Besides their use as a testbed, corpora here are beginning to make significant contributions to the actual capabilities of systems. This is because of the increasing availability of bilingual corpora (e.g. in the proceedings of institutions that are legally required to be available in multiple languages: literary classics and best-sellers may also often be available in this form, though the standards of equivalence in translation may be laxer). These enable a self-organizing approach to supplement the traditional knowledge-based one. Speech processing in general is benefiting more and more from the development of its own speech corpora. These compilations, though, are still quite distinct from text corpora, being extended acoustic analyses of speech wave-forms; hence they are not intrinsically based on character or word analysis as is a text corpus. Furthermore, the amount of language represented in such a corpus is several orders of magnitude smaller than a typical corpus of texts. As such, they are beyond the scope of this paper.

Nevertheless, text corpora derived from speech will be increasingly of use to the developers of speech processors. Such corpora will be literal transcriptions of spoken language. Part of their advantage for speech analysts lies in their potentially much greater scale: textual-type annotation is much quicker and easier to add the type of annotations required for a speech corpus. But this does not in itself advance the capabilities of speech processing. Where it will help is in making it possible to identify some of the higher-level regularities correlating with parameters of a given type of speech (e.g. proneness to hesitation in certain contexts, or preferred sentence structure) : these in turn can be used to tune the speech processor.

9.5 Types of Use

By way of summary, we can discern two major classes of use for corpora amidst all this potential variety: the corpus as a large-scale, but heavily-diluted source of data, which new techniques are enabling us to sift; and the corpus as a testbed, composed of representative if mostly undifferentiated material, good for testing or training an automatic device under development.

Although we have described a corpus as undifferentiated, this is meant only in the sense

that the precise content of the component texts is not to the point. For the corpus to satisfy as a useful testbed, it must be highly calibrated, with as many of the external parameters of its constituent texts stated as possible. The different users can then either discover the linguistic correlates of these external parameters; or else use them to judge which parts of the corpus are appropriate as a test for their processing device.

Select Bibliography

- Aarts, J. & W. Meijs (eds.) (1986) *Corpus Linguistics II. Recent Advances in the Use of Computer Corpora in English Language Research*. Amsterdam: Rodopi.
- Atkins B. T. (1987) 'Semantic ID tags: corpus evidence for dictionary senses', in: *The Uses of Large Text Databases: Proceedings of the 3rd Annual Conference of the UW Centre for the New Oxford English Dictionary*. Canada: University of Waterloo.
- Biber, D. (1989) 'A typology of English texts', in: *Linguistics* 27.
- Church K., W. Gale, P. Hanks, & D. M. Hindle (1990) 'Using statistics in lexical analysis', in: Uri Zernik (ed.) *Lexical Acquisition: Using On-line Resources to Build a Lexicon*. Lawrence Erlbaum.
- Clear, J. (1988) 'Trawling the language: monitor corpora', in: M. Snell-Hornby (ed.) *ZüriLEX '86 Proceedings: Papers read at the EURALEX International Congress*. Tübingen: Francke Verlag.
- Clear (forthcoming) 'Corpus sampling', Proceedings of 11th ICAME Conference, Berlin 1990.
- Engwall, G. (forthcoming) 'Chance or choice: criteria for corpus selection' in B. T. Atkins & A. Zampolli (eds.) *Computational Approaches to the Lexicon*, OUP.
- Francis, W. Nelson (1979) 'Problems of assembling and computerizing large corpora', in: H. Bergenholz & B. Shäder (eds.) *Empirische Textwissenschaft: Aufbau und Auswertung von Text-Corpora*. Königstein: Scriptor Verlag.
- Francis, W. Nelson (1980) 'A tagged corpus—problems and prospects', in: S. Greenbaum, G. Leech and J. Svartvik (eds.) *Studies in English Linguistics, for Randolph Quirk*. London and New York: Longman.
- Garside R., G. Leech, & G. Sampson (1987) *The Computational Analysis of English*. London: Longman.
- Halliday, M. (1966) 'Lexis as a linguistic level', in: C. Bazell, J. Catford, M. Halliday and R. Robins (eds.) *In Memory of J. R. Firth*. London: Longman.
- Johansson, S. (1980) 'The LOB corpus of British English texts: presentation and comments', in: *ALLC Journal* 1.
- Johansson, S. (ed.) (1982) *Computer Corpora in English Language Research*. Bergen: Norwegian Computing Centre for the Humanities.
- Källgren, G. (1990) 'The first million is the hardest to get', in: *Proceedings of the 13th International Conference on Computational Linguistics (COLING)*. Helsinki: Association for Computational Linguistics.
- Klavans, J., & E. Tzoukermann (1990) 'Linking bilingual corpora and machine readable dictionaries with the BICORD system', in: *Proceedings of the 6th Annual Conference of the Centre for the New OED and Electronic Text Research*. Canada: University of Waterloo.
- Kucera, H., & W. N. Francis (1967) *Computational Analysis of Present-Day American English*. Providence, RI: Brown University Press.
- Kyto, M. (1989) 'Introduction to the use of the Helsinki corpus of English texts: diachronic and dialectal', in: M. Ljung (ed.) *Proceedings of the Computer Conference in Stockholm*.
- Kyto, M., O. Ihalainen, & M. Rissanen (1988) *Corpus Linguistics, Hard and Soft: Proceedings of the 8th International Conference on English Language Research on Computerized Corpora*. Amsterdam: Rodopi.

- Leitner, G. (1990) 'Corpus design—problems and suggested solutions', working paper in: *ICE Newsletter* 7, May 1990. University College London: International Corpus of English.
- Leitner, G., & U. Schäfer (1989) 'Reflections on corpus linguistics—the 9th ICAME conference in Birmingham, England', in: *CCE Newsletter* 3(1).
- Liberman, M. (1989) 'Text on tap: the ACL/DCI', in: *Proceedings of DARPA Speech and Natural Language Workshop*. New York: Morgan Kaufman.
- Meijs, W. (ed.) (1987) *Corpus Linguistics and Beyond: Proceedings of the 7th International Conference on English Language Research on Computerized Corpora*. Amsterdam: Rodopi.
- Oostdijk, N. (1988) 'A corpus for studying linguistic variation', in: *ICAME Journal* no.12.
- Oostdijk, N. (1988) 'A corpus linguistic approach to linguistic variation', in: *Literary and Linguistic Computing* 3.
- Renouf, A. (1987) 'Corpus Development', in: J. Sinclair (ed.) *Looking Up*. London: Collins.
- Quirk, R., & J. Svartvik (1979) 'A corpus of modern English', in: H. Bergenholz & B. Shäder (eds.) *Empirische Textwissenschaft: Aufbau und Auswertung von Text-Corpora*. Königstein: Scriptor Verlag.
- Shastri, S. V. (1988) 'The Kolhapur corpus and work done on its basis so far', in: *ICAME Journal* no.12.
- Sinclair, J. (1982) 'Relection on computer corpora in English language research', in: S. Johansson (ed.) *Computer Corpora in English Language Research*. Bergen: Norwegian Computing Centre for the Humanities.
- Sinclair, J. (ed.) (1987) *Looking Up*. London: Collins.
- Sinclair, J. (1989) 'Corpus creation', in: Candlin & McNamara (eds.) *Language, Learning and Community*. Sydney: NCELTR Macquairie University.
- Svartvik, J. (1990) *The London-Lund Corpus of Spoken English: Description and Research*. Lund Studies in English, 82. Lund University Press.
- The Uses of Large Text Databases: Proceedings of the 3rd Annual Conference of the UW Centre for the New Oxford English Dictionary*. (1987). Canada: University of Waterloo.
- Warwick, S., & J. Hajic (1990) 'Searching on tagged corpora: linguistically motivated concordance analysis', in: *Electronic Text Research: Proceedings of the Conference of the Centre for the New OED and Electronic Text Research*. Canada: University of Waterloo.
- Woods, A., P. Fletcher, & A. Hughes (1986) *Statistics in Language Studies*. Cambridge: CUP.