

Applications of software in the compilation of corpora

Raymond Hickey, Department of English, University of Munich

Abstract

An attempt is made here to sketch some of the applications to which corpus processing software can be put in the compilation of corpora. The emphasis is on the one hand on the automation of many standard processes, such as text collation and the provision of header information for each file of a corpus, while on the other hand the additional possibilities offered by dedicated corpus software are described. Among the latter special emphasis is put on the transfer of textual data to a database environment for further processing. Further matters such as the use of special fonts for older stages of English and the option of organising the text files of one's corpus for potential users in advance are also discussed.

0 Introduction

Given the nature of the contributions to this volume, the present author thought it appropriate to discuss the uses to which corpus processing software could be put in the compilation and distribution of corpora, especially ones with a diachronic orientation.

Assuming that the compilers of a corpus have reached basic agreement on what periods are to be covered and what texts are to be included, software can be used gainfully from this point onwards. To illustrate possible applications the software system *Lexa* developed by the present author will be used for the ensuing discussion. This programme suite and the three volumes of documentation are available from the Norwegian Computer Centre for the Humanities in Bergen, Norway. Each of the following sections is intended to illustrate a typical situation in which software is useful in the preparatory stage of corpus building. The list is not exhaustive but it does cover the main areas of concern in this phase of text collection and organisation.

1 Text collation

It is safe to assume that more than one individual will be involved in the compilation of a text corpus. Texts will either be scanned or keyed in directly. In either case it is more the exception than the rule to find that a text turns up error-free in the computer. This banal fact increases the status of the individual who is responsible for text correction. Again it is commonplace for more than one version of a text to exist in some intermediary stage of compilation. Sooner or later in such a situation doubts arise as to whether a particular version of a text is the more accurate or the better corrected. The need arises quite quickly for a reliable means of comparing two versions of a single text. Of course the time and date stamp of a file on the operating system level will tell you which of two is the more recent, but age is not a guarantee for correctness.

To resolve this dilemma a programme has been included in the *Lexa* suite which will

linguistic straightjacket as it imposes the grammatical classification scheme of the compilers on the user. Seeing that there is tagging software available, many compilers may now prefer to leave this work to the corpus users, or to some sub-group, such as researchers in another university who would be prepared to carry out this task. As universities have to economise on resources, tagging by the compilers is likely to become less likely in future, especially as partial tagging is not viewed as a sensible course of action. You either tag completely or not at all. If you decide to do so you may bind your capacities in a manner which you come to regret later.

This would appear at least to be the case for major projects like the Helsinki corpus. With the arrival of smaller more specialised corpora, tagging may become feasible particularly if it is directly connected with the research interests of the corpus compilers.

4 Using Cocoa headers

Independent of the question of whether to tag or not to tag, compilers of a corpus should consider whether it would be of avail to future users to include some information on the nature of the texts before distributing these. This decision has fortunately been made in favour of supplying such information by the compilers of the Helsinki corpus. The format they have chosen for the inclusion of text-relevant information is what is commonly known as the Cocoa header. Note that header information is placed at the top of a file and has nothing to do with grammatical classifications included in the body of a text.

Parameters of the Cocoa header

1: <B = 'name of text file'>	2: <Q = 'text identifier'>
3: <N = 'name of text'>	4: <A = 'author'>
5: <C = 'part of corpus'>	6: <O = 'date of original'>
7: <M = 'date of manuscript'>	8: <K = 'contemporaneity'>
9: <D = 'dialect'>	10: <V = 'verse' or 'prose'>
11: <T = 'text type'>	12: <G = 'relation to foreign original'>
13: <F = 'foreign original'>	14: <W = 'relation to spoken language'>
15: <X = 'sex of author'>	16: <Y = 'age of author'>
17: <H = 'social rank of author'>	18: <U = 'audience description'>
19: <E = 'participant relation'>	20: <J = 'interaction'>
21: <I = 'setting'>	22: <Z = 'prototypical text category'>
23: <S = 'sample'>	24: <P = 'page'>
25: <L = 'line'>	26: <R = 'record'>

Although providing headers is not comparable to the task of tagging a corpus, it nonetheless requires an additional amount of work to specify the values for these parameters for the texts of a corpus. The advantages, however, are considerable.

With the *Lexa* suite the contents of a Cocoa header can be accessed by the information retrieval software. This is realised as follows: a programme (called *Cocoa*) extracts the header information from any set of input files and deposits this in a database. Then with the database manager of the suite (called *DbStat*) one can load the database just created and impose a filter on it. By this is meant that only those records remain visible which match a

have the shapes of the Old and Middle English characters) in any set of input texts at those points where it encounters an escape sequence, e.g. it inserts the *yogh* symbol when it hits on '+g'. The conversion is reversible so that texts can be restored to their original form if desired. The numerical values of the redefined characters with Old and Middle English shapes are as following.

Escape sequence	Actual character	Letter name	ASCII numerical value for redefinition by Lexa programme	Make Symbols
+a	æ	'ash',l.c.	145	+A
	Æ	'Ash',u.c.	146	+d
		'eth',l.c.	253	+D
		'Eth',u.c.	252	+g
				'yogh',l.c.
			243	+G
242	+t	"	'thorn',l.c.	245
			'Thorn',u.c.	244
+tt	ø	'crossed thorn'	248	+TT
		'crossed Thorn'	246	—
+e	É	'e caudata'	144	+l
		'pound sign'	156	œ

Taking a typical text such as *Beowulf* and loading it with one's text editor or word processor leads to one being presented with a text which is convenient for computer manipulation but hardly readable to the Old English scholar.

Beginning of Beowulf with 'escape-sequence' coding

```
[ ] [\BEOWULF\ ] ] <R 1> Hw+at. We Gardena in geardagum, +teodcyninga, +trym
gefrunon, hu +da +a+telingas ellen fremedon. <R 4> Oft Scyld Scefing [{scea+tena{]}
+treatum, monegum m+ag+tum, meodosetla ofteah, egsode eorlas. <R 6> Sy+d+dan +arest
[{\wear+d{]} feasceaft funden, he +t+as frofre gebad, weox under wolcnum,
weor+dmyndum +tah, o+d+t+at him +aghwylc +tara ymbsittendra ofer hronrade hyran
scolde, gomban gyldan. <R 11> +t+at w+as god cyning. <R 12> +d+am eafera w+as +after
cenned, geong in geardum, +tone god sende folce to frofre; fyren+dearfe ongeat +te hie +ar
drugon [{aldorlease{]} lange hwile. <R 16> Him +t+as liffrea, wuldres wealdend,
woroldare forgeaf; Beowulf w+as breme bl+ad wide [{sprang{]} , Scyldes eafera
Scedelandum in.
```

If one takes this text, however, and runs it through the programme *Make Symbols* which is supplied in the *Lexa* suite then a number of substitutions are made and certain high ASCII characters are inserted where escape sequences were found in an input text. Now under the important assumption that (i) you are using a computer with a colour monitor, typically a VGA video adapter based system, and (ii) that you have loaded the supplied Old / Middle English font of the *Lexa* suite then the stretch of text printed above should now look like the following.

- Oxford: Oxford University Computing Service.
- Johansson, Stig 1986. *The tagged LOB Corpus: User's manual*. Bergen: Norwegian Computing Centre for the Humanities.
- Johansson, Stig and Anna-Brita Stenström (eds.). 1991. *English computer corpora: Selected papers and research guide*. Berlin: Mouton de Gruyter.
- Kytö, Merja, Ossi Ihalainen, and Matti Rissanen (eds.) 1988. *Corpus linguistics hard and soft*. Amsterdam: Rodopi.
- Kytö, Merja 1991. *Manual to the diachronic part of the Helsinki corpus of English texts*. Helsinki: Department of English.
- Kytö, Merja and Matti Rissanen 1988. The Helsinki Corpus of English Texts: Classifying and coding the diachronic part. In *Corpus linguistics*, ed. by M.Kytö et al. 169-180. Amsterdam: Rodopi.