

Cambridge University Press
978-0-521-83851-1 - Corpus Linguistics: Method, Theory and Practice
Tony Mcenery and Andrew Hardie
Frontmatter
[More information](#)

Corpus Linguistics

Corpus linguistics is the study of language data on a large scale – the computer-aided analysis of very extensive collections of transcribed utterances or written texts. This textbook outlines the basic methods of corpus linguistics, explains how the discipline of corpus linguistics developed, and surveys the major approaches to the use of corpus data. It uses a broad range of examples to show how corpus data has led to methodological and theoretical innovation in linguistics in general. Clear and detailed explanations lay out the key issues of method and theory in contemporary corpus linguistics. A structured and coherent narrative links the historical development of the field to current topics in ‘mainstream’ linguistics. Practical activities and questions for discussion at the end of each chapter encourage students to test their understanding of what they have read and an extensive glossary provides easy access to definitions of all technical terms used in the text.

TONY MCENERY is Professor of English Language and Linguistics at Lancaster University.

ANDREW HARDIE is Lecturer in Corpus Linguistics at Lancaster University.

Cambridge University Press
978-0-521-83851-1 - Corpus Linguistics: Method, Theory and Practice
Tony Mcenery and Andrew Hardie
Frontmatter
[More information](#)

Cambridge University Press
978-0-521-83851-1 - Corpus Linguistics: Method, Theory and Practice
Tony Mcenery and Andrew Hardie
Frontmatter
[More information](#)

CAMBRIDGE TEXTBOOKS IN LINGUISTICS

General editors: P. AUSTIN, J. BRESNAN, B. COMRIE, S. CRAIN,
W. DRESSLER, C. EWEN, R. LASS, D. LIGHTFOOT, K. RICE,
I. ROBERTS, S. ROMAINÉ, N. V. SMITH

Corpus Linguistics:
Method, Theory and Practice

In this series:

S. C. LEVINSON *Pragmatics*
G. BROWN and G. YULE *Discourse Analysis*
R. HUDDLESTON *Introduction to the Grammar of English*
R. LASS *Phonology*
B. COMRIE *Tense*
W. KLEIN *Second Language Acquisition*
A. J. WOODS, P. FLETCHER and A. HUGHES *Statistics in Language Studies*
D. A. CRUSE *Lexical Semantics*
A. RADFORD *Transformational Grammar*
M. GARMAN *Psycholinguistics*
G. G. CORBETT *Gender*
H. J. GIEGERICH *English Phonology*
R. CANN *Formal Semantics*
J. LAVER *Principles of Phonetics*
F. R. PALMER *Grammatical Roles and Relations*
M. A. JONES *Foundations of French Syntax*
A. RADFORD *Syntactic Theory and the Structure of English: A Minimalist Approach*
R. D. VAN VALIN, JR, and R. J. LAPOLLA *Syntax: Structure, Meaning and Function*
A. DURANTI *Linguistic Anthropology*
A. CRUTTENDEN *Intonation* Second edition
J. K. CHAMBERS and P. TRUDGILL *Dialectology* Second edition
C. LYONS *Definiteness*
R. KAGER *Optimality Theory*
J. A. HOLM *An Introduction to Pidgins and Creoles*
G. G. CORBETT *Number*
C. J. EWEN and H. VAN DER HULST *The Phonological Structure of Words*
F. R. PALMER *Mood and Modality* Second edition
B. J. BLAKE *Case* Second edition
E. GUSSMAN *Phonology: Analysis and Theory*
M. YIP *Tone*
W. CROFT *Typology and Universals* Second edition
F. COULMAS *Writing Systems: An Introduction to their Linguistic Analysis*
P. J. HOPPER and E. C. TRAUGOTT *Grammaticalization* Second edition
L. WHITE *Second Language Acquisition and Universal Grammar*
I. PLAG *Word-Formation in English*
W. CROFT and A. CRUSE *Cognitive Linguistics*
A. SIEWIERSKA *Person*
A. RADFORD *Minimalist Syntax: Exploring the Structure of English*
D. BÜRING *Binding Theory*
M. BUTT *Theories of Case*
N. HORNSTEIN, J. NUÑES and K. GROHMANN *Understanding Minimalism*
B. C. LUST *Child Language: Acquisition and Growth*
G. G. CORBETT *Agreement*
J. C. L. INGRAM *Neurolinguistics: An Introduction to Spoken Language Processing and its Disorders*
J. CLACKSON *Indo-European Linguistics: An Introduction*
M. ARIEL *Pragmatics and Grammar*
R. CANN, R. KEMPSON and E. GREGOROMICHELAKI *Semantics: An Introduction to Meaning in Language*
Y. MATRAS *Language Contact*
D. BIBER and S. CONRAD *Register, Genre and Style*
L. JEFFRIES and D. MCINTYRE *Stylistics*
R. HUDSON *An Introduction to Word Grammar*
M. L. MURPHY *Lexical Meaning*
J. M. MEISEL *First and Second Language Acquisition*
T. MCENERY and A. HARDIE *Corpus Linguistics: Method, Theory and Practice*

Cambridge University Press
978-0-521-83851-1 - Corpus Linguistics: Method, Theory and Practice
Tony Mcenery and Andrew Hardie
Frontmatter
[More information](#)

Corpus Linguistics

Method, Theory and Practice

TONY McENERY AND
ANDREW HARDIE

Lancaster University



Cambridge University Press
978-0-521-83851-1 - Corpus Linguistics: Method, Theory and Practice
Tony Mcenery and Andrew Hardie
Frontmatter
[More information](#)

CAMBRIDGE
UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom

Cambridge University Press is part of the University of Cambridge.

It furthers the University’s mission by disseminating knowledge in the pursuit of education, learning and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/9780521547369

© Tony McEnery and Andrew Hardie 2012

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

First published 2012

3rd printing 2014

Printed in the United Kingdom by Clays, St Ives plc.

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

Library of Congress Cataloguing in Publication data

McEnery, Tony, 1964–

Corpus linguistics : method, theory and practice / Tony McEnery, Andrew Hardie.

p. cm. – (Cambridge textbooks in linguistics)

Includes index.

ISBN 978-0-521-83851-1 (hardback)

1. Corpora (Linguistics) I. Hardie, Andrew. II. Title.

P128.C68M38 2011

410.1’88 – dc23 2011026519

ISBN 978-0-521-83851-1 Hardback

ISBN 978-0-521-54736-9 Paperback

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this publication, and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

Contents

<i>List of figures</i>	<i>page</i> x
<i>List of tables</i>	xi
<i>Acknowledgements</i>	xii
<i>Preface</i>	xiii
1 What is corpus linguistics?	1
1.1 Introduction	1
1.2 Mode of communication	3
1.3 Corpus-based versus corpus-driven linguistics	5
1.4 Data collection regimes	6
1.5 Annotated versus unannotated corpora	13
1.6 Total accountability versus data selection	14
1.7 Monolingual versus multilingual corpora	18
1.8 Summary	21
Further reading	21
Practical activities	22
Questions for discussion	23
2 Accessing and analysing corpus data	25
2.1 Introduction	25
2.2 Are corpora the answer to all research questions in linguistics?	27
2.3 Corpus annotation	29
2.4 Introducing concordances	35
2.5 A historical overview of corpus analysis tools	37
2.6 Statistics in corpus linguistics	48
2.7 Summary	53
Further reading	54
Practical activities	55
Questions for discussion	55
3 The web, laws and ethics	57
3.1 Introduction	57
3.2 The web and legal issues	57
3.3 Ethical issues	60
3.4 Summary	69
Further reading	69
Practical activity	70
Questions for discussion	70

viii	Contents
4	English Corpus Linguistics 71
4.1	Introduction 71
4.2	University College London (UCL) 74
4.3	Lancaster University 76
4.4	University of Birmingham 79
4.5	Université Catholique de Louvain 81
4.6	University of Nottingham 84
4.7	Northern Arizona University and the USA 88
4.8	Summary 91
	Further reading 91
	Practical activities 92
	Questions for discussion 92
5	Corpus-based studies of synchronic and diachronic variation 94
5.1	Introduction 94
5.2	Diachronic change from Old English to Modern English 94
5.3	Diachronic variation in contemporary Modern English 96
5.4	The multi-dimensional approach to variation 104
5.5	Corpora and variationist sociolinguistics 115
5.6	Summary 118
	Further reading 119
	Practical activities 119
	Questions for discussion 120
6	Neo-Firthian corpus linguistics 122
6.1	Introduction 122
6.2	Collocation 122
6.3	Discourse 133
6.4	Semantic prosody and semantic preference 135
6.5	Lexis and grammar 142
6.6	Corpus-as-theory versus corpus-as-method 147
6.7	Summary: Sinclair's contribution to corpus linguistics 162
	Further reading 164
	Practical activities 164
	Questions for discussion 165
7	Corpus methods and functionalist linguistics 167
7.1	Introduction 167
7.2	Functionalism in linguistics: a brief overview 168
7.3	Corpus-based research from a functionalist perspective 171
7.4	Corpora and typology 176
7.5	Corpora and cognitive approaches to linguistics 179
7.6	Corpora in the analysis of metaphor 185
7.7	Summary 188
	Further reading 189
	Practical activities 189
	Questions for discussion 191

	Contents	ix
8 The convergence of corpus linguistics, psycholinguistics and functionalist linguistics	192	
8.1 Introduction	192	
8.2 Corpus methods and psycholinguistics	193	
8.3 The convergence of neo-Firthian corpus linguistics and functionalist linguistics	210	
8.4 Summary	221	
Further reading	222	
Practical activities	223	
Questions for discussion	224	
9 Conclusion	225	
9.1 Introduction	225	
9.2 The story of corpus linguistics, from past to future	225	
9.3 Revisiting old friends: computational linguistics	227	
9.4 The textually mediated world: the humanities and social sciences	230	
9.5 The challenge ahead: integrating corpora with new methods in linguistics	233	
9.6 The final word	236	
<i>Glossary</i>	238	
<i>Notes</i>	254	
<i>References</i>	259	
<i>Index</i>	292	

Figures

2.1	A conversation between June and Jonathan (BNC file KCT, utterances 357–365)	<i>page</i> 30
2.2	An extract of a sample concordance of the particle 了 from the LCMC	35
2.3	An extract of a concordance of words ending in <i>-ness</i> from BE06 (Baker 2009)	36
3.1	A thief revealed	63
4.1	Spoken transcript from the CANCODE Corpus, from McCarthy and Carter (2001: 52–3)	86
5.1	Biber’s Dimension 2, <i>Narrative versus Non-Narrative Concerns</i> (from Biber 1988)	107
5.2	A fragment of a feature tree for English	114

Tables

1.1	The LOB Corpus Sampling Frame (after Hofland and Johansson 1982: 2)	page 10
1.2	A hypothetical corpus	23
2.1	Metadata stored about two speakers, June and Jonathan, in BNC file KCT	29
4.1	Corpus annotation research at Lancaster University in the 1980s and 1990s	78
4.2	Features whose frequency differentiates speech and writing (all examples taken from Leech (1998: 11–13))	88
5.1	The Brown Corpus sampling frame	97
5.2	The ‘Brown Family’ of corpora	99
6.1	Three word n-grams (‘lexical bundles’) beginning with <i>cheese</i> with a frequency of ten or more in the written section of the BNC	124
6.2	The top ten collocates of <i>cheese</i> in the BNC calculated using different statistical measures (intra-sentential collocates within a span of +/–3 only)	128
6.3	The top ten collocates of <i>cheese</i> in the BNC calculated using the same statistical measure (log-likelihood) but different spans	128
6.4	The top ten colligates of <i>cheese</i> in the written BNC calculated using a word-based and a tag-based approach (statistic: log-likelihood; span: +/–3, intra-sentential collocates only)	131
6.5	Hunston and Francis’ analysis of interlocking patterns in a sample sentence	144

Acknowledgements

This book could not have been written without the aid of many colleagues whose generous assistance we gratefully acknowledge. Svenja Adolphs, Karin Aijmer, Mark Davies, Costas Gabrielatos, Geoffrey Leech, Neil Millar and Richard Zhonghua Xiao all read part or all of the manuscript and offered useful criticisms and suggestions (although, as should go without saying, responsibility for the final text rests solely with us). We would also like to thank Eivind Torgerson for a helpful discussion regarding the similarity of corpus linguistics and socio-linguistics, and Ghada Mohamed for long talks on the topics of cluster analysis and the notion of an exhaustive linguistic feature tree (both of which fed into the discussion in Chapter 5). Others, too numerous to mention, among our colleagues and students pointed us towards relevant literature or gave us valuable sneak pre-views of books and research papers in press which we would not otherwise have been able to discuss here. We are thankful to them all. We are also grateful for the professionalism of the editorial staff at Cambridge University Press.

Preface

The title of this book is *Corpus Linguistics: Method, Theory and Practice*. As that title may suggest, it is about how corpus linguistics has developed and is employed as a methodology; the major theoretical issues that corpus linguists contend with today; and the problems that researchers using corpora must grapple with in practice, both within linguistics and across disciplines.

This captures in a nutshell, we like to think, what makes this book different from the many excellent introductory textbooks on corpus linguistics that have appeared since the mid-1990s. Our purpose in writing this book is not to introduce the very basics of procedures in corpus linguistics – we do not outline any step-by-step instructions describing how to go about investigating a corpus, and though we do occasionally give some example corpus analyses, we do not go into the details of how to deal with concordances, collocations, keywords and other common outputs from corpus tools. Other books have covered this ground comprehensively (e.g. Biber *et al.* 1998; Hunston 2002; Adolphs 2006; McEnery *et al.* 2006; Hoffmann *et al.* 2008), and we do not see any need to duplicate these accounts.

Instead, our aim in this textbook is to introduce, explain and in some cases problematise the most fundamental conceptual issues underlying the use of corpora, as well as reviewing what we see as the major trends of research using corpora to date. In the earlier part of the book, we will discuss corpus linguistics as a discipline, exploring high-level issues of practice that are of concern to corpus linguists, such as: features of corpus construction such as the notions of balance and representativeness; the development and exploitation of corpus tools; issues of copyright law and of research ethics; the role and limits of corpus annotation; the role of quantitative analysis; and so on. We will also review the research priorities and main contributions of a number of different schools and centres of corpus linguistics. In doing so, we will explore both their methodological apparatus and, where appropriate, the contributions to theory that they have sought to make. However, in later chapters we will move beyond the boundaries of corpus linguistics *per se* to consider the role that corpus methods now have across a range of types of linguistic investigation – including studies of language variation, language change, functional-cognitive linguistic theory, and psycholinguistics – and the various issues raised by the use of corpora in such enterprises. Accordingly we will argue that, although ‘corpus linguistics’ has clearly long had a separate existence, there is a very great degree of convergence between

corpus linguistics and these other aspects of linguistics. Corpus techniques tend no longer to be the preserve of a clearly delimited field of specialists, but rather have become a critical resource across linguistics as a whole (and beyond). Thus, we might argue that the future of the field is in ‘corpus methods in linguistics’ rather than ‘corpus linguistics’ standing independently.

This view, it must be noted, is not universally held. In particular, some scholars of the neo-Firthian school of corpus linguistics disagree entirely with it, as we will explain in our discussion of that tradition in Chapter 6. When covering matters such as this, which do not attract full consensus, we had two choices as authors. The first was to attempt to maintain neutrality, and to write the text without letting our opinions and theoretical and methodological preferences colour the account. The second was to acknowledge our perspective, and associated biases, frankly and explicitly, and to provide an account in which we explain and justify our views as best we can. It is this second approach which we have adopted. Our discussion of the different traditions of corpus research, of the wider use of corpora in linguistics and beyond, and of the future directions we see as desirable, thus amount to what might be called a position statement for our ‘version’ of corpus linguistics (to borrow the memorable phrasing of Teubert 2005). We make no pretence to neutrality, and our discussions of traditions in corpus linguistics other than our own may indeed characterise a given school of thought in terms that the researchers within that tradition would not necessarily agree with. The most notable case of this is that our discussion of neo-Firthian corpus linguistics is very much informed by our stance as non-neo-Firthian researchers. We urge the reader who is interested in understanding the full picture with regard to these debates to refer to accounts in which the scholars we discuss have dealt with these issues from their own perspective, with their different sets of preconceptions, opinions and biases; where appropriate we list such accounts in the suggestions for further reading provided in each chapter. In the case of neo-Firthian theory, for instance, we would not hesitate to recommend Tognini-Bonelli (2001) or Teubert and Čermáková (2004) as clear and readable presentations of the ‘other side’ of the argument that we make here in Chapters 6 and 8.

Our recommendations for further reading cover both the primary and the secondary literature, as appropriate. We have also included some study questions at the end of each chapter. These are divided into two groups. For readers who are interested in thinking further about some of the issues and problems that we cover, we suggest a number of *questions for discussion*. Alongside that, we also provide some *practical activities*. These activities provide practice in some of the key procedures that are needed in actually ‘doing’ corpus linguistics. All the practical activities assume that you have access to a reasonably large corpus and some corpus analysis software, but do not require any specific corpus. So if you have a copy of a suitable corpus on your computer, you can work with that data using any one of several different corpus tools to complete the exercise (some of which are available for free on the Internet). If not, then a number

of large, standard corpora – in several languages – can be accessed and analysed via the World Wide Web, using online interfaces which either are openly accessible, or offer freely available sign-in accounts. There are too many of these to list here (but see our discussion in Chapter 2). For English, we recommend BNCweb (<http://bncweb.lancs.ac.uk/bncwebSignup>) or the Brigham Young University online interface (<http://corpus.byu.edu>). While different software tools have different capabilities, the activities in this book are based on the core functions that all corpus tools make available.

The centrality of the Internet to the practice of corpus linguistics today means that it was necessary to make reference to websites and web services at various points in the book. We have tried to keep the number of web addresses in the text to a minimum, in the light of how variable these can be over relatively short periods of time. It is inevitable, however, that some of these web addresses will become outdated over time. For this reason, we have established a companion website for this book, where updated links will be made available where necessary. The website will also contain other supplementary material, including suggested answers to the study questions in each chapter. Importantly, this site also contains a large number of additional notes that were simply too numerous to include in the book. Accordingly, we suggest that you check the website after you read each chapter. The companion website address is www.cambridge.org/mcenery-hardie.