

Corpus Stylistics

Speech, writing and thought
presentation in a corpus of English
writing

Elena Semino and Mick Short

 **Routledge**
Taylor & Francis Group
LONDON AND NEW YORK

**Also available as a printed book
see title verso for ISBN details**

Corpus Stylistics

This book combines stylistic analysis with corpus linguistics in order to provide an innovative account of the phenomenon of speech, writing and thought presentation – commonly referred to as ‘speech reporting’ or ‘discourse presentation’.

This new account is based on an extensive analysis of a quarter-of-a-million word electronic collection of written narrative texts, including both fiction and non-fiction. The book includes detailed discussions of:

- The construction of a corpus of late twentieth-century written British narratives, taken from fiction, newspaper news reports and (auto)biographies.
- The development of a manual annotation system for speech, writing and thought presentation and its application to the corpus.
- The findings of a quantitative and qualitative analysis of the forms and functions of speech, writing and thought presentation in the three genres represented in the corpus.
- The findings of the analysis of a range of specific phenomena, including hypothetical speech, writing and thought presentation, embedded speech, writing and thought presentation, and ambiguities in speech, writing and thought presentation.
- Two case studies concentrating on specific texts from the corpus.

Corpus Stylistics shows how stylistics, and text/discourse analysis more generally, can benefit from the use of a corpus methodology. The authors’ innovative approach results in a more reliable and comprehensive categorization of the forms of speech, writing and thought presentation than has been suggested so far. This book will be essential reading for linguists interested in the areas of stylistics and corpus linguistics.

Elena Semino is Senior Lecturer in the Department of Linguistics and Modern English Language at Lancaster University. She is the author of *Language and World Creation in Poems and Other Texts* (1997), and co-editor (with Jonathan Culpeper) of *Cognitive Stylistics: Language and Cognition in Text Analysis* (2002). **Mick Short** is Professor of English Language and Literature at Lancaster University. He has written *Exploring the Language of Poems, Plays and Prose* (1996) and (with Geoffrey Leech) *Style in Fiction* (1981). He founded the Poetics and Linguistics Association, and was the founding editor of its international journal, *Language and Literature*.

Routledge advances in corpus linguistics

Edited by Anthony McEnery

Lancaster University, UK

and

Michael Hoey

Liverpool University, UK.

Corpus-based linguistics is a dynamic area of linguistic research. The series aims to reflect the diversity of approaches to the subject, and thus to provide a forum for debate and detailed discussion of the various ways of building, exploiting and theorizing about the use of corpora in language studies.

1 Swearing in English

Anthony McEnery

2 Antonymy

A corpus-based perspective

Steven Jones

3 Modelling Variation in Spoken and Written English

David Y. W. Lee

4 The Linguistics of Political Argument

The spin-doctor and the wolf-pack at the White House

Alan Partington

5 Corpus Stylistics

Speech, writing and thought presentation in a corpus of English writing

Elena Semino and Mick Short

Corpus Stylistics

Speech, writing and thought
presentation in a corpus of English
writing

Elena Semino and Mick Short

First published 2004
by Routledge
11 New Fetter Lane, London EC4P 4EE

Simultaneously published in the USA and Canada
by Routledge
29 West 35th Street, New York, NY 10001

Routledge is an imprint of the Taylor & Francis Group

This edition published in the Taylor & Francis e-Library, 2004.

© 2004 Elena Semino and Mick Short

All rights reserved. No part of this book may be reprinted or reproduced or utilized in any form or by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying and recording, or in any information storage or retrieval system, without permission in writing from the publishers.

British Library Cataloguing in Publication Data

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

Library of Congress Cataloging in Publication Data

A catalog record for this book has been requested

ISBN 0-203-49407-5 Master e-book ISBN

ISBN 0-203-57142-8 (Adobe eReader Format)

ISBN 0-415-28669-7 (Print Edition)

Contents

<i>List of figures</i>	viii
<i>List of tables</i>	ix
<i>Acknowledgements</i>	xi
1 Introduction: a corpus-based approach to the study of discourse presentation in written narratives	1
<i>1.1 Introduction</i>	<i>1</i>
<i>1.2 Why a corpus-based approach?</i>	<i>4</i>
<i>1.3 The Leech and Short (1981) model</i>	<i>9</i>
<i>1.4 Other corpus based approaches to speech, writing and thought presentation</i>	<i>16</i>
<i>1.5 The structure of this book</i>	<i>17</i>
2 Methodology: the construction and annotation of the corpus	19
<i>2.1 The corpus</i>	<i>19</i>
<i>2.2 The annotation system</i>	<i>26</i>
<i>2.3 Concluding remarks</i>	<i>39</i>
3 A revised model of speech, writing and thought presentation	42
<i>3.1 New categories and a new presentational scale</i>	<i>42</i>
<i>3.2 New sub-categories</i>	<i>52</i>
<i>3.3 An overview of speech, writing and thought presentation in the corpus</i>	<i>57</i>
<i>3.4 Concluding remarks</i>	<i>64</i>

4	Speech presentation in the corpus: a quantitative and qualitative analysis	66
	4.1 <i>Introduction</i>	66
	4.2 <i>The speech presentation categories in the corpus</i>	66
	4.3 <i>Concluding remarks</i>	96
5	Writing presentation in the corpus: a quantitative and qualitative analysis	98
	5.1 <i>Introduction</i>	98
	5.2 <i>The writing presentation categories in the corpus</i>	98
	5.3 <i>Concluding remarks</i>	111
6	Thought presentation in the corpus: a quantitative and qualitative analysis	114
	6.1 <i>Introduction</i>	114
	6.2 <i>The pure thought presentation categories in the corpus</i>	116
	6.3 <i>Inferred thought presentation in the corpus</i>	135
	6.4 <i>Concluding remarks</i>	147
	6.5 <i>An overview of our findings on the major speech, writing and thought presentation categories</i>	149
7	Specific phenomena in speech, writing and thought presentation	153
	7.1 <i>Quotation phenomena</i>	153
	7.2 <i>Hypothetical speech, writing and thought presentation</i>	159
	7.3 <i>Embedded speech, writing and thought presentation</i>	171
	7.4 <i>Ambiguity in speech, writing and thought presentation</i>	182
	7.5 <i>Concluding remarks</i>	198
8	Case studies of specific texts from the corpus	201
	8.1 <i>Introduction</i>	201
	8.2 <i>Is the medium the message? The presentation of conversations with the dead in Joyful Voices by Doris Stokes</i>	202
	8.3 <i>Discourse presentation in newspaper reports of a 'PC Bible' story</i>	210
9	Conclusion	222
	9.1 <i>Our findings and the corpus approach</i>	222
	9.2 <i>Areas where further research is needed</i>	227

<i>Appendix 1</i>	List of texts sampled	232
<i>Appendix 2</i>	The SW&TP tagset	235
<i>Appendix 3</i>	Alphabetical list of reporting verbs for Indirect Speech presentation	237
<i>Appendix 4</i>	Alphabetical list of reporting verbs for Direct Speech presentation	239
<i>Appendix 5</i>	Alphabetical list of reporting verbs for Indirect Writing presentation	242
<i>Appendix 6</i>	Alphabetical list of reporting verbs for Direct Writing presentation	243
<i>Appendix 7</i>	Alphabetical list of reporting verbs for Direct Thought presentation	244
<i>Appendix 8</i>	Alphabetical list of reporting verbs for Indirect Thought presentation	245
<i>Bibliography</i>		246
<i>Index</i>		251

Figures

1.1	The speech presentation scale	11
1.2	The 'norm' on the speech presentation scale	13
1.3	The thought presentation scale	14
1.4	The speech and thought presentation scales and their respective 'norms'	15
3.1	The speech, writing and thought presentation scales	49
8.1	Diagrammatic representation of alternative sets of beliefs deriving from the assumption that there is a 'spirit world'	204
8.2	Diagrammatic representation of alternative sets of beliefs deriving from the assumption that there is no 'spirit world'	205

Tables

3.1	Numbers of occurrences of speech, writing, thought and other tags in the corpus	59
3.2	Percentages of speech, writing, thought and other tags out of all tags in the corpus	59
3.3	Percentages of words included under the speech, writing, thought and other tags out of all words in the corpus	59
3.4	Percentages of speech, writing, thought and other tags out of all tags in the six sub-sections of the corpus	64
4.1	Numbers of occurrences of the speech presentation categories in the corpus	67
4.2	Mean word length of the speech presentation categories in the corpus	68
4.3	Numbers of occurrences of NRSA and NRSAp in the corpus	74
4.4	Numbers of occurrences of DS and FDS tags in the corpus	91
5.1	Numbers of occurrences of the writing presentation categories in the corpus	100
5.2	Mean word length of the writing presentation categories in the corpus	101
5.3	Numbers of occurrences of DW and FDW tags in the corpus	112
6.1	Numbers of occurrences of the thought presentation categories in the corpus	115
6.2	Numbers of pure (i.e. non-inferred) thought presentation categories in the corpus	117
6.3	Numbers of occurrences of DT and FDT tags in the corpus	121
6.4	Mean length of the thought presentation categories in the corpus	122

6.5	Numbers of occurrences of inferred thought presentation categories in the corpus	137
6.6	Relative proportions of inferred thought presentation categories in the biography and autobiography sections of the corpus	139
7.1	The eight most frequent 'q' tags in the corpus	156
7.2	Hypothetical SW&TP tags in the corpus	168
7.3	Hypothetical SW&TP tags in the three genres included in the corpus	169
7.4	Occurrences of embedded speech presentation categories in the corpus	176
7.5	Occurrences of embedded writing presentation categories in the corpus	177
7.6	Occurrences of embedded thought presentation categories in the corpus	178
7.7	The 25 most frequent portmanteau tags	184
8.1	'PC Bible' stories	211
8.2	DS, DW and 'q' in the 'PC Bible' articles	217

Acknowledgements

The research presented in this book was supported by grants from the Faculty of Social Sciences at Lancaster University and from the *Humanities Research Board* of the British Academy (grant BA LRG M-AN2314/AON3489).

We are grateful to the following publishers for permission to draw from parts of our previously published papers:

Pearson Education for permission to draw from: Short, M., Semino, E. and Culpeper, J. (1996) 'Using a corpus for stylistics research: speech and thought presentation', in Thomas, J. and Short, M. (eds) *Using Corpora in Language Research*, London: Longman, pp. 110–31; C. Winter for permission to draw from: Short, M., Wynne, M. and Semino, E. (1999) 'Reading reports: discourse presentation in a corpus of narratives, with special reference to news reports', in Diller, H.-J. and Stratmann, E. O.-J. (eds) *English Via Various Media*, Heidelberg: Winter, pp. 39–65; Sage Publications for permission to draw from Short, M., Semino, E. and Wynne, M. (2002) 'Revisiting the notion of faithfulness in discourse presentation using a corpus approach', *Language and Literature*, 11, 4, 325–55 © Sage Publications Ltd, 2002; Elsevier Science for permission to draw from Semino, E., Short, M. and Culpeper, J. (1997) 'Using a corpus to test a model of speech and thought presentation', *Poetics*, 25, 17–43; Peter Lang for permission to draw from Short, M. (2003) 'A corpus-based approach to speech, thought and writing presentation', in Wilson, A., Rayson, P. and McEnery, T. (eds) *Corpus Linguistics by the Lune: A Festschrift for Geoffrey Leech*, Frankfurt/Main: Peter Lang. We draw from a small section of: Wynne, M., Short, M. and Semino, E. (1998) 'A corpus-based investigation of speech, thought and writing presentation in English narrative texts', in Renouf, A. (ed.) *Explorations in Corpus Linguistics*, Amsterdam: Rodopi, 231–45. We also draw from parts of the following paper: Semino, E., Short, M. and Wynne, M. (1999) 'Hypothetical words and thoughts in contemporary British narratives', *Narrative* 7, 3: 307–34. Copyright 1999 by the Ohio State University Press. All rights reserved.

We are grateful to a number of individuals for contributing in significant ways to the project which led to this book. Ruth Allen and Markus

Guadagnin worked with us during the pilot phase of the project. Jonathan Culpeper was involved with our initial analysis of the corpus, and was co-author for two of our joint papers. Martin Wynne was our Research Assistant during the main phase of the project and collaborated with us on most of the papers that resulted from the project.

Several other students and colleagues assisted us at various points during our research: Mike Dodgson, Salah El-Hassan, Eleni Gogorosi, Reiko Ikeo, Scott Piao, Itzumi Tanaka and Richard Xiao. Our two Research Assistants on a closely related project, Daniel McIntyre and John Heywood, have provided us invaluable assistance and insight. We are also grateful to our long-suffering corpus linguistics colleagues at Lancaster University, who kindly answered our many questions: Paul Baker, Geoff Leech, Tony McEnery and Nick Smith. Graeme Hughes and Damien Cashman patiently got us out of several technical difficulties. As series co-editor, Michael Hoey gave us invaluable feedback and advice on an earlier version of the manuscript.

We have greatly benefited from the feedback we have received from the audiences at many conferences where we have presented papers based on the project described in this book, and from referees' and editors' comments on our articles and book chapters. Finally, we are grateful to the many generations of students who have forced us to clarify our thinking by asking questions, or simply looking puzzled during our classes.

1 Introduction

A corpus-based approach to the study of discourse presentation in written narratives

1.1 Introduction

We hope that this book will be of interest to at least two different kinds of linguists: (i) textlinguists (e.g. stylisticians and critical discourse analysts) who are involved in the analysis of discourse presentation in written and spoken language, and (ii) corpus linguists or other linguists who are interested in developing dedicated electronic corpora to elucidate textual phenomena. As we try to take both of these main readerships into account, we may, to some degree, tell one readership what it already knows. We apologize in advance if we sometimes do this, and we will try to keep such descriptions to a minimum. Nonetheless, we think it helpful to try to draw the textlinguistic and corpus traditions closer together through this specific study.

Our book describes the research on discourse presentation in written narratives we have been involved in since 1994, and which is still ongoing.¹ This work has involved the systematic and detailed annotation of a corpus of written fictional and non-fictional narratives for speech, writing and thought presentation categories, in order to throw light on discourse presentation theory and on how patterns of discourse presentation vary in three different written narrative genres (fiction, news reports and (auto)biographies).

Since 1996 we have published seven articles and book chapters on our work.² However, because these articles are spread through different books and journals, it is difficult for scholars to access the reports of the work we have undertaken. This volume, which draws from parts of these articles but also contains new material, is a summation of our work to date – work which aims to offer insights in relation to the study of discourse presentation in texts and to what is a relative innovative methodology for textlinguists. We will also use this book to consolidate what has been for us a constantly developing method of textual annotation and theory building. Because our research project has evolved over time, our articles to date have some descriptive and annotational inconsistencies among them. We have gradually changed some of the terms and annota-

2 Introduction

tions we have used as we have come to grips with new discourse presentation phenomena in our data. These inconsistencies may well have been confusing for those who have read more than one of our articles, and this volume provides an opportunity to explain the changes we have made and our reasons for making them, and to arrive at a reasonably stable set of descriptive terms and annotations for further research. We do not, of course, assume that our work to date is the end of the story in descriptive, annotational, analytical or theoretical terms.³ We hope that others might be interested in applying the analytical methods we have developed to yet other spoken and written genres/text types,⁴ to see how well our approach works for these other genres and how the patterns of discourse presentation in these genres compare with those we have analysed.

Before we proceed further, it will be helpful if we make some points about our use of terminology in this book. We have used the term ‘discourse’ in the discussion above for two reasons. First, we sometimes need a general, and briefer, term to refer to what we otherwise call ‘speech, writing and thought presentation’ (SW&TP).⁵ We will strive to use the term ‘discourse presentation’ only in this general, overarching sense. Our second reason for using the term was that we wanted to connect our work to that of other scholars who have written about the way in which the discourse of others is presented, and who often use the term ‘discourse presentation’ for this enterprise. However, we are conscious of the fact that the term ‘discourse’ is often used vaguely and/or with somewhat different meanings by different scholars. We have pointed out before (Short *et al.* 2002) that one of the dangers of the term ‘discourse presentation’ is that, if it is used as an elegant variant of the more specific terms ‘speech presentation’, ‘writing presentation’ and ‘thought presentation’, it is possible to move seamlessly from the discussion of one mode of presentation to another without making the change clear to oneself, or to others. This in turn can lead to mis-analyses and a less accurate understanding of the phenomena under investigation. We believe that, although there are commonalities among speech presentation, writing presentation and/or thought presentation, there are also important differences which are unhelpfully hidden if the general term ‘discourse presentation’ is used as an alternative for these more specific, mode-related terms and concepts. Hence, when discussing specific discourse presentation phenomena, we will strive to use the more specific terms and not to use the general term as a substitute for them.

The other term which we have already made considerable use of is ‘presentation’. We use this term as a default, rather than ‘report’ or ‘representation’ (which are often used as default terms by other linguists), because we are specifically interested in how the discourse of others (or the speaker/writer on some previous occasion) is *presented*. This is what textual annotation and analysis can most sensibly be used for (and

explains why stylisticians tend to use this term). We prefer not to use the term 'report', which is often used as a default by grammarians (e.g. Huddleston and Pullum 2002: 1023–30; Quirk *et al.* 1985: 1020–33) and other linguists who are part of a tradition where examples are invented when discussing discourse presentation. This is because the term 'report' suggests an unproblematic relationship between the discourse presentation and the anterior discourse which is being presented. Tannen (1989), among others, has shown that an assumption of faithful report for direct speech presentation in casual conversation is unrealistic (yet interestingly she uses the term 'report' even when undermining this assumption). However, we do not want to use the term 'representation' as a default either, as this tends to be used by linguists (e.g. critical discourse analysts like Caldas-Coulthard 1994 and Fairclough 1988) who want to concentrate mainly on distortions and misrepresentations in the reporting of anterior discourses. 'Presentation' is thus helpfully neutral for the discussion of speech, writing and thought presentation in a corpus of written texts where, for the most part, we do not, in any case, have easy access to the anterior speech, writing or thought being presented. We discuss this issue of terminology in more detail in Short *et al.* (2002).⁶

Many studies have proposed models of the forms and functions of discourse presentation in a range of text-types (e.g. Bally 1912a, 1912b; Banfield 1982; Collins 2001; Fairclough 1988; Fludernik 1993; Fowler 1986; McHale 1978; Pascal 1977; Tannen 1989; Thompson 1994, 1996; Volosinov 1973; Waugh 1995; see also papers in Coulmas 1986 and Lucy 1993). The original motivation for our corpus-based study of discourse presentation, however, was to test how well the particular model of speech and thought presentation outlined in Leech and Short (1981: Ch. 10) worked on written text types other than the novel. The Leech and Short model was developed specifically to account for the range of speech and thought presentation forms and their effects in novels written in English. We wanted to test this model, not only because one of us has a rather obvious personal interest in it, but also because (i) it is still the most analytically specific account of speech and thought presentation to date, and (ii) it has been influential and widely used by other textlinguists.

Many analysts of prose fiction, including Fludernik (1993: 283–316, *passim*) and Simpson (1993: 21–30), have discussed the Leech and Short approach. Person (1999: 28–37) and Toolan (2001: 136–40) also include discussions of some of our more recent work referred to above. A number of studies have also applied the Leech and Short approach to non-literary texts. McKenzie (1987) uses Leech and Short to analyse how free indirect speech was used to circumvent a ban on direct quotation of the ANC in a booklet by South African students, and Roeh and Nir (1990) use it in the analysis of Israeli radio broadcasts. Thompson's (1996) account of the dimensions of choice available to speakers or writers when reporting the language of others also draws on the Leech and Short model, which

he describes as ‘comprehensive in its coverage’ and ‘[t]he most fully developed’ of the various approaches to speech and thought presentation (Thompson 1996: 504).

1.2 Why a corpus-based approach?

The Leech and Short model, like all theoretical models in stylistics up to that point, was developed through the use of scholarly intuition, based on extensive personal reading experience, which was in turn exemplified and tested through the analysis of examples chosen from previous reading. The model was also designed to account specifically for speech and thought presentation in fictional texts (indeed, most of the discourse presentation work by stylisticians and narratologists has concentrated on fiction). Hence it was difficult to know how generalizable the model was to other text-types, or how descriptively adequate it was when ‘tested to destruction’ on texts (including fictional texts) in a way that could not avoid inconvenient or borderline cases. It was for this reason that we decided to develop and annotate a dedicated corpus to test out the model.

We should also point out that some of the non-corpus work on discourse presentation which has already been completed has been based on the accumulation of very large numbers of examples accrued from previous reading. Specific mention should be made here of the monumental work of the narratologists Cohn (1978) and Fludernik (1993). We have benefited considerably from these two very insightful works. Cohn grounded her analysis of what we would call thought presentation through the accumulation of a manually collected corpus of examples:

Equipped with these basic abstractions [of narrative theory] I could then travel around in narrative literature, selecting works and passages in works that would best display the entire spectrum of possibilities, while in turn allowing these works themselves to reveal unforeseen hues.

(Cohn 1978: v)

Cohn’s motivation is not unlike ours, except that we want to compare discourse presentation across text types, including narrative fiction, and want to be much more explicit about our criteria for text selection, as well as being more explicit and systematic in our analysis of the texts in our corpus. Cohn was writing before computers could be used to store and interrogate large corpora of texts, of course, and we could well imagine that if she were beginning her work now, she might also want to make use of an electronic corpus, as we have.

Fludernik’s (1993) study of what she calls free indirect discourse is even more impressive in terms of the wide range of textual examples she uses

to illustrate the points she wants to make. We have learned much from her work but, as with Cohn's study, we were concerned that her relatively informal analytical approach might mean that important factors in the study of discourse presentation would be missed. In her research, Fludernik specifically considered the possibility of a corpus-based approach, and the quantification that comes with it, but rejected this option (i) because she did not want to restrict herself to the literature of just one language, nation, period, etc., which she thought a corpus-based approach would prevent, and (ii) because she believed that a corpus and its associated annotation would have created serious methodological problems, in the sense that she thinks it would have been necessary to 'institute arbitrary definitions of the relevant categories' (Fludernik 1993: 9):

Such arbitrariness would necessarily have resulted in an erosion of the actual usefulness of the statistical data, since one would have had either to decide on larger categories that include marginal and ambiguous phenomena, or to indulge in a proliferation of subcategories and intermediary categories which would have rendered the statistics next to useless for interpretation. From previous experience with statistical research (Fludernik 1982) I have also acquired a profound distrust of the methodological relevance of statistical data. Statistics typically take individual occurrences of certain phenomena out of context. Since the present study attempts to document the crucial importance of context for the purpose of the even preliminary establishment of basic categories, a statistical approach would from the outset have vitiated one of the major aims of the project. These remarks are, however, not meant to discredit statistical research in itself. On the contrary, I would welcome a series of statistical analyses that might help to corroborate, modify or refute some of the theses I am here proposing.

(Fludernik 1993: 9)

We have quoted from Fludernik at length because we have effectively tried to do what she decided to avoid, namely to use a set of categories and subcategories to analyse the textual extracts in our corpus comprehensively and systematically. Consequently, we certainly recognize some of the problems she points to, though we think that the annotation difficulties have not been as damaging as she thought they would be. Indeed, we would claim that forcing ourselves to be as clear and precise as possible about our annotations has helped us to isolate, and come to terms with, phenomena we may not otherwise even have noticed. Similarly, we believe that forcing ourselves to account for ambiguity and marginal phenomena in our annotations has helped us to understand more exactly how the speech, writing and thought presentation scales operate, and what factors are at work in producing ambiguity on those scales. Because we take this

explicit analytical approach, we are able to provide some of the statistical information which, at the end of the above quotation, Fludernik says that she would welcome.

We very much agree with Fludernik that statistical analysis has limitations as well as advantages, and this is why we present both quantitative and qualitative analysis in this book. We do not think that the one precludes the other (though doing both does increase the workload still further, as, from experience, we are very well aware). Indeed, we would want to argue that both forms of analysis are needed, and work best when used interdependently. Although Fludernik decided not to adopt a corpus-based and quantitative approach (the experience of the dissertation she refers to as Fludernik 1982 was clearly salutary!), she makes a point of saying that she is not antipathetic to such work. She is very open to the fact that all approaches have advantages and disadvantages, and that we can all learn from different approaches to the same phenomenon. This tolerant and inclusive attitude is in contrast to the attacks on corpus linguistics by some other linguists, which we allude to briefly below.

It was natural for us to move to a corpus-based approach as we work in a department which has members who have been involved in corpus construction and annotation for some years, and who could easily be called upon for advice and help. The Lancaster–Oslo/Bergen (LOB) corpus was one of the early modern linguistic corpora to be developed; Lancaster is the ‘home’ of the British National Corpus (BNC), for which Lancaster did much of the work, and our colleagues are involved in the building and exploitation of other corpora too. However, not all linguists are sympathetic to a corpus-based approach, and so we will take a little space here to explore some of the pros and cons in the use of electronic corpora, to help explain our decision to develop our corpus and to use ‘corpus stylistics’ as the main title of this book.

The first point that we would like to make is that although this book, and much of our current work, involves the use of a corpus-based approach in stylistics, we do not think that this approach should supplant other work within our field. Rather, our decision to use a corpus-based approach was because it was the best tool we could find to carry out the particular kind of investigation we had in mind. In order to see how adequate the Leech and Short model was, and what kind of modifications it might require, we needed to test the model on a number of different text-types, with enough samples of each text-type to be reasonably sure of our findings. This led to the idea of a representative corpus. We also needed to force ourselves not just to concentrate on convenient text- or intuition-based examples. This led to the idea of developing a method of systematic and replicable textual annotation which would be used comprehensively. Finally, we needed to be able to sort our annotations easily, in order to observe patterns of various kinds in our data. This need led naturally to the idea of using an electronic tagged corpus, and software that would

enable us to do what we wanted (we chose Mike Scott's Wordsmith package for this purpose).

The fact that we are currently involved in corpus-based work, and the quantification that it entails, does not mean that we have stopped doing the qualitative textual analysis that is at the heart of the field of stylistic analysis. We are still involved in this sort of work, and will continue to do it (indeed this book includes some qualitative work on particular texts in our corpus; see Chapter 8 in particular). We will continue to use our intuition in arriving at theories, interpretations of texts and so on, and we will not give up our interest in investigating informant reactions to texts in order to compare them with stylistic analyses or stylistic theories – or indeed any other kind of work we, or other stylisticians, typically engage in. We think that all these different approaches have a useful role to play in helping us (i) to understand how readers interact with, and understand, particular texts and (ii) to arrive at general theories of textual understanding, textual response and style. We would be unhappy if the work we report was regarded as a competitor for other forms of enquiry in stylistics, rather than as merely another (very useful) approach to add to the analytical armoury of the stylistics enterprise. There is already some interesting work which insightfully combines detailed qualitative work on particular texts with corpus-based analysis. Stubbs (1996: 81–100) uses such a combination to show how Baden-Powell, the founder of the Boy Scout movement, uses the same lexical items in very different (and sexist) ways in his last messages to the Girl Guides and the Boy Scouts, and Louw (1997) uses corpus-based work to show, for example, how what he calls the 'semantic prosody' of the word 'utterly' is used by Philip Larkin to induce feelings of threat at the end of 'First Sight', his poem about newborn lambs:

They could not grasp it if they knew,
What so soon will wake and grow
Utterly unlike the snow.

To some it may seem strange that we need to point out at all what we have just said in the last paragraph. However, corpus linguistics has, in recent years, been part of the sort of 'turf war' that breaks out from time to time in most academic disciplines. These rather heated debates have mainly concerned the use of large generalized corpora (e.g. the Brown Corpus, the LOB Corpus, the Bank of English and the BNC), which were set up to investigate empirically the lexical and grammatical characteristics of English and other languages. We do not have the space to enter into this debate here, and in any case our corpus is not a generalized corpus of the sort over which the arguments have raged, but a much smaller affair, set up with a much more specific set of goals. For those interested in the debate over generalized corpora, McEnery and Wilson (1996: 1–18)

provide a useful account of the history of the relationship between corpus linguistics and ‘mainstream’ linguistics. For more recent contributions to the ‘corpus linguistics wars’, see Borsley and Ingham’s (2002) attack on corpus linguistics from a ‘theoretical linguistics’ standpoint, to which Stubbs (2002) responds, and Widdowson’s (2000) critique from an ‘applied linguistics’ perspective.⁷

Not surprisingly, these sorts of debates are often characterized by misunderstandings and caricatures of others’ positions. Academic turf wars tend to generate more heat than understanding. We prefer to take the more cooperative and inclusive view of Biber *et al.* (1998: 7–8), who argue that corpus-based analysis ‘should be seen as a complementary approach to more traditional approaches, rather than as the single correct approach’, and of Fillmore (1992), who says:

I don’t think there can be any corpora, however large, that contain information about all of the areas of the English lexicon and grammar that I want to explore . . . [but] every corpus I have had the chance to examine, however small, has taught me facts I couldn’t imagine finding out any other way. My conclusion is that the two types of linguists need one another.

(Fillmore 1992: 35, quoted in McEnery and Wilson 1996: 25)

Stylisticians have always occupied a fairly peripheral position in the panoply of linguistic description and theorizing (it is the literary critics who have worried rather more about us occupying part of their territory), and so we do not feel particularly personally affected by the antagonistic debates between the corpus linguists and other linguists. That said, it is difficult to believe that the study of the linguistic performance seen in texts, at least, can be adequately conducted without the use of the corpus-based approach, in addition to other approaches.

We hope that we have already made it clear that we are interested in combining corpus-based techniques with more intuition-based approaches. Our corpus work could not have been as successful as it has if considerable prior intuition-based work on discourse presentation was not available for us to test and develop, as we hope our discussion of Cohn’s and Fludernik’s work above, and that of Leech and Short (1981) in 1.3 below, makes clear. Indeed, in terms of general research design, it makes sense to move first from intuition-based and genre-specific work to corpus-based work on a range of reasonably close genres. This is what we have done so far, and report in this book. We compare fictional texts with news report and biography and autobiography. If the model being examined (and revised) operates successfully in these areas, the next steps will be to expand the work to cover other written genres and also to test the methodology out on spoken data (see note 3). To extend the work to naturally occurring spoken interaction is bound to be more difficult, if only

because of the need to take account of the turn-taking phenomena and normal non-fluency typical of spoken discourse, and we would argue that it would be difficult to undertake such work successfully without the prior work on the more 'orderly' medium of writing which we have been carrying out.

1.3 The Leech and Short (1981) model

As we pointed out in 1.1 above, our annotation work is based on the scales of speech and thought presentation outlined in Leech and Short (1981: Ch. 10; see also Short 1996: Ch. 10). Indeed, a primary aim of our corpus work was to see whether the Leech and Short model, which had been developed to account for the meanings and effects of the speech and thought presentation categories in literary prose fiction, could be applied sensibly, systematically and with insight to non-literary and non-fictional narrative modes. The Leech and Short model was the first to distinguish systematically between the presentation of speech and the presentation of thought in the novel. It also suggested, as some other scholars did (e.g. Cohn 1978 and McHale 1978; see also Fludernik 1993: 283–4), that the discourse presentation scales are not an assemblage of hard-edged, discrete categories, but continua, rather like that seen in the colour spectrum. The speech and thought presentation scales had the same categories and in the same order along the scales, but Leech and Short pointed out that some of the categories had different effects on the different scales (in particular, free indirect thought had effects which were often opposite to those for free indirect speech, and the direct and free direct forms had different effects in speech, as opposed to thought, presentation). The Leech and Short account also suggested a new category (the narrative report of speech acts, and its equivalent on the thought presentation scale) and the re-positioning of the free direct category on the scales. Instead of being positioned between the free indirect and direct categories, as assumed by scholars previously, Leech and Short proposed that the free direct category (free direct speech, free direct thought) was at one extreme end of the scales, 'beyond' the direct forms (direct speech and direct thought). Since 1981 this ordering appears to have been generally accepted by most scholars (see, for example, Fludernik 1993: 289–315, Simpson 1993: 21–30 and Toolan 2001: 116–40, but contrast Person 1999: 19–32).

For reasons of clarity, we will first concentrate on the speech presentation scale. It had been traditionally assumed that direct speech (DS) and indirect speech (IS) were distinguished not just in terms of their formal linguistic features, but also in terms of whether the words and grammatical structures of the original utterance were presented, as well as its propositional form. Leech and Short, building on the work of earlier stylisticians, saw the entire speech presentation scale (which was already

known to have more categories than just IS and DS) as being ordered in relation both to the linguistic features involved and also to the number of faithfulness claims with respect to the original that the speaker's/writer's choice of speech presentation category involved. The speech presentation category distinctions given below are ordered on a scale which relates to the amount of 'involvement' of (i) the original speaker in the anterior discourse and (ii) the person in the posterior discourse presenting what was said in the anterior discourse (bold typeface is used to indicate the specific stretch of text that exemplifies each category). Because the Leech and Short descriptive system was construed mainly in relation to the novel, the 'original speakers' were characters and the reporters were narrators (hence the use of 'N' for 'Narrative' and 'Narration' in the abbreviations below):

- (N) = Narration – no speech presentation involved (hence the bracketing of the symbol here)
 e.g. **He looked straight at her.**
- NRSA** = Narrative Report of Speech Acts
 e.g. He looked straight at her and **told her about his imminent return.** She was pleased.
- IS** = Indirect Speech
 e.g. He looked straight at her and told her **that he would definitely return the following day.** She was pleased.
- FIS** = Free Indirect Speech
 e.g. He looked straight at her. **He would definitely come back tomorrow!** She was pleased.
- DS** = Direct Speech
 e.g. He looked straight at her and said **'I'll definitely come back tomorrow!'**.
- FDS** = Free Direct Speech
 e.g. He looked straight at her. **'I'll definitely come back tomorrow!'** She was pleased.

Narration sentences (presenting states, events and actions in the fictional world) are not strictly part of the speech presentation scale and so 'N' is placed in brackets above. It is usually included in the presentation of such scales because NRSA is linked closely with N, being the presentation of speech as action. The speech presentation categories can be distinguished to a large degree in linguistic terms. Readers will be familiar with the orthographic, syntactic and deictic distinctions between IS and DS, so we will not reiterate them here. For Leech and Short, FDS must obligatorily contain the direct string, but need not contain either the reporting clause or the punctuation surrounding the direct string. It is because they regard these features as being provided by the narrator/reporter in written presentations of speech that they argue that FDS should be at one extreme of

the scale. In its most extreme form, it presents the words of the character/original speaker with no apparent 'interference' from the narrator/reporter.

NRSA, unlike IS, prototypically has only one clause, with the 'speech report' verb often followed by a noun phrase or a prepositional phrase indicating the topic of the speech presented. Because this kind of presentation is more minimal than the propositional form associated with indirect strings, NRSA is placed between N and IS on the above scale. Not surprisingly, NRSA is prototypically used for summarizing, and for providing background speech information to contextualize fuller speech presentation forms.

Free indirect speech (FIS) is a form between IS and DS because it shares linguistic features associated prototypically with both the IS and DS forms. Typically, it will not have the quotation marks associated with DS and often does not have the reporting clause associated with IS. It may contain some deictic features (in the widest sense of the term) which are appropriate for DS and, at the same time, others which are appropriate for IS (cf. 'tomorrow' vs 'he' in the above FIS examples). In contrast to previous scholars, Leech and Short argued that no particular linguistic features were criterial for FIS to occur. All you needed was a mix of the sorts of features normally associated with DS and IS. Previous scholars had assumed that third-person pronouns and backshift of tense compared with the associated DS form were criterial for FIS. But Leech and Short pointed out that these features were effectively neutralized in first-person narrations and present-tense narrations respectively, and so could not be criterial in all cases.

This issue of criteriality is an important one for us. In tagging our corpus we used formal criteria to distinguish categories as much as we were able because they are the most reliable criteria to apply consistently. However, the application of formal features – 'rules' – does not always yield an analysis which works, and indeed it is possible to find cases where, formally, a particular sentence (or sentence part) could belong to more than one category, and only the application of contextual considerations can yield a satisfactory assignment, if one can be found at all (see 9.1.1 for further comments on these issues).

The speech and thought presentation scales are usually represented as being ordered along a horizontal axis, with NRSA in the left-most speech presentation position, adjacent to (N) and the free direct category in the right-most position. Hence the speech presentation continuum is usually represented visually as in Figure 1.1.

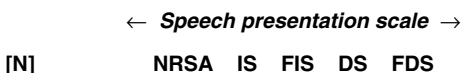


Figure 1.1 The speech presentation scale.

At the extreme ends of the speech presentation scale shown in Figure 1.1 we get (1) narration, where no speech presentation is involved at all and (2) (free) direct speech (i.e. FDS and DS together), where it is assumed canonically by readers that the direct string reports exactly the words and structures used by the character to say whatever they said in the ‘anterior’ discourse. The narrative report of speech acts (NRSA)⁸ category was thought by Leech and Short to be the ‘hinge’ between speech presentation and narration (speech acts are both speech and actions). In NRSA the speech act value of the utterance presented is indicated, often with a specification of the topic of the speech act, but no more elaboration of what was said in the anterior discourse is made. Thus, in marked contrast to (F)DS, this ‘summarizing’ nature of NRSA displays a fairly loose connection with both *what* was said (its propositional content) and *how* it was said (the words and structures used to utter the relevant propositional content). In the Leech and Short account of speech presentation in the novel (where the narrator presents what characters have said, or say), indirect speech (IS) displays a greater ‘contribution’ from the character in the novel than NRSA because it makes a weightier claim to be faithful to the original. As we indicated above, NRSA tells us the speech act value of what was said, plus a specification (sometimes optional) of the topic of the speech act. IS does this and, *in addition*, presents the propositional content of what was said. The use of (F)DS normally brings one further faithfulness claim: in addition to presenting the speech act value and the propositional content of the utterance, it provides the words and grammatical structures claimed to have been used to utter the propositional content and associated speech act. This extra faithfulness claim brings with it associated effects of vividness and dramatization. Hence an (F)DS representation of some speech in a novel, for example, feels foregrounded, vivid and immediate as compared with an IS version.

The functional notion of increasing degrees of faithfulness to an original, as one moves from left to right on the speech presentation continuum, helps to explain why it is that we have such a full panoply of presentational forms when we write. We should remember, though, that it is open to writers to misuse the canonical forms, for example by using the DS form but not using the words and structures uttered in some original, in order to mislead or rhetorically affect readers (see Chapter 8 and Short *et al.* 2002). Moreover, as Leech and Short pointed out, fiction is unusual in discourse presentation terms. Most discourse presentation involves an anterior discourse which is re-presented in the posterior, reporting discourse. However, this is not normally true in fictions where there is no actual anterior speech to be presented. The whole story, including the account of ‘what was said earlier’ is fictional, and we merely pretend ‘conventionally’ that the conversation ‘reported’ took place in the world of the fiction.

Most of what we have said so far is well known to stylisticians, but the

free indirect speech (FIS) category in particular may be new to others. It is a crucial category for stylisticians because it is often associated with ironic effects when it is used to present character speech in fiction. In quantitative terms, the proportion of FIS in our corpus is small compared with the other major speech presentation categories. However, its equivalent on the thought presentation scale, free indirect thought (FIT), is used very extensively in the novel, and is the most frequent of Leech and Short's thought presentation categories in both the fiction section of our corpus and the corpus more generally. Effectively, FIS is a 'mix' of the deictic and other features associated with IS on the one hand and DS on the other, and as a consequence is ambiguous with respect to the 'words and structures' faithfulness claim. It is often difficult to know, for particular words, whether they 'belong' to the character or the narrator/reporter. If we take the FIS example used above (He would definitely come back tomorrow!), it is clear that if a narrator or reporter is presenting what someone else previously said, the third-person pronoun and the backshifted modal verb would normally be assumed to 'belong' to that narrator/reporter because the expressions are deictically inappropriate for the original speaker, who would normally use 'I' to refer to himself, for example. Because 'come back' and 'tomorrow' are deictically proximal, it will often be assumed that they must 'belong' to the original speaker, particularly when, as in the examples above (where the DS and FDS forms can be compared with the FIS one), the context is set up to encourage that assumption. However, if the narrator/reporter happens to be presenting what was originally said on the same day and in the same place as the original utterance, then 'come back' and 'tomorrow' will be deictically appropriate both for the original utterance and its posterior presentation. Similarly, the 'exclamatory tone' suggested by the exclamation mark could be attributed either to the original utterance or its posterior presentation.

Because FIS is a 'deictic mix' of the words of the original and its presentation by someone else, Leech and Short (1981), who had argued that the norm for speech presentation is DS, went on to suggest that FIS is perceived by readers as distancing them from what the character said (often with attendant effects of irony), since its choice constitutes a movement away from the DS norm (see Figure 1.2) towards the narrator/reporter end of the scale.

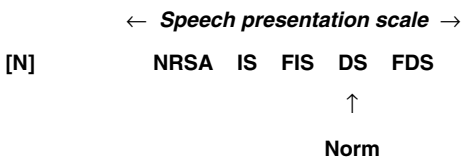


Figure 1.2 The 'norm' on the speech presentation scale.

In other words, FIS is the nearest category to DS in which readers feel that the narrator ‘interposes’ him- or herself between the words of the character and the reader. We will return to the notion of ‘norms’ for speech and thought presentation below and in Chapters 4 and 6.

Compared with previous accounts, then, the Leech and Short (1981) model, besides being more explicit, also established the NRSA category and reorganized the categories into an order which related to the faithfulness claims. This enabled a more orderly and principled account of the presentational effects obtained when a writer uses one presentation category rather than another. The definitions of categories for speech presentation were partly on functional grounds (the faithfulness claims), partly on linguistic grounds (made as explicitly as possible) and partly on contextual grounds (for example, sometimes sentences can be formally ambiguous between narration and free indirect speech but unambiguous when interpreted in context).

As we said above, Leech and Short also distinguished in a clear way for the first time between speech presentation and thought presentation. They set up a separate scale of thought presentation, with categories parallel to those on the speech presentation scale, and defined in analogous ways (see Figure 1.3).

Below we give prototypical examples for the thought presentation categories to match those we provided earlier for speech presentation (note that the free indirect and the free direct examples can be formally identical to their speech presentation equivalents, but, when situated in appropriate co-text, it would be clear contextually that they were presenting thought, not speech):

(N) = Narration – no thought presentation involved (hence the bracketing of the symbol here)
e.g. **He looked straight at her.**

NRTA = Narrative Report of Thought Acts
e.g. He looked straight at her and **thought about his imminent return**. She remained unaware of his plan until the following day.

IT = Indirect Thought
e.g. He looked straight at her and decided **that he would definitely return the following day**. She remained unaware of his plan until the following day.

FIT = Free Indirect Thought
e.g. He looked straight at her. **He would definitely come back tomorrow!** She remained unaware of his plan until the following day.