# Corpora in Applied Linguistics

Corpus linguistics has revolutionised the world of language study and is an essential component of work in applied linguistics. This book, now in its second edition, provides a thorough introduction to all the key research issues in corpus linguistics, from the point of view of applied linguistics. The field has progressed a great deal since the first edition, so this edition has been completely rewritten to reflect these advances, whilst still maintaining the emphasis on hands-on corpus research of the first edition. It includes chapters on qualitative and quantitative research, applications in language teaching, discourse studies, and beyond. It also includes an extensive discussion of the place of corpus linguistics in linguistic theory, and provides numerous detailed examples of corpus studies throughout. Providing an accessible but thorough grounding to the fascinating, fast-moving field of corpus linguistics, this book is essential reading for the student and the researcher alike.

Susan Hunston is Professor of English Language at the University of Birmingham, UK. Her recent publications include *Corpus Approaches to Evaluation* (Routledge, 2011) and *Interdisciplinary Research Discourse* (with Paul Thompson, Routledge, 2019). She is a former Chair of the British Association for Applied Linguistics.

## THE CAMBRIDGE APPLIED LINGUISTICS SERIES

The authority on cutting-edge Applied Linguistics research

Series Editors       2007–present: Carol A. Chapelle and Susan Hunston
                             1988–2007: Michael H. Long and Jack C. Richards

For a complete list of titles please visit: www.cambridge.org

*Recent titles in this series*

**Corpora in Applied Linguistics**
2nd edition
*Susan Hunston*

**The Language of Mental Illness**
Corpus Linguistics and the Construction
of Mental Illness in the Press
*Hazel Price*

**Mobile Assisted Language Learning**
Concepts, Contexts and Challenges
*Glenn Stockwell*

**Research Genres Across Languages**
Multilingual Communication Online
*Carmen Pérez-Llantada*

**Validity Argument in Language Testing**
Case Studies of Validation Research
*Edited by Carol A. Chapelle and Erik Voss*

**Doing English Grammar**
Theory, Description and Practice
*Roger Berry*

**Learner Corpus Research Meets Second
Language Acquisition**
*Bert Le Bruyn and Magali Paquot*

**Second Language Speech Fluency**
From Research to Practice
*Parvaneh Tavakoli and Clare Wright*

**Ontologies of English**
Conceptualising the Language for Learning,
Teaching, and Assessment
*Edited by Christopher J. Hall and Rachel
Wicaksono*

**Task-Based Language Teaching**
Theory and Practice
*Rod Ellis, Peter Skehan, Shaofeng Li,
Natsuko Shintani and Craig Lambert*

**Feedback in Second Language Writing**
Contexts and Issues
2nd edition
*Edited by* Ken Hyland and Fiona Hyland

**Language and Television Series**
A Linguistic Approach to TV Dialogue
*Monika Bednarek*

**Intelligibility, Oral Communication, and the
Teaching of Pronunciation**
*John M. Levis*

**Multilingual Education**
Between Language Learning and
Translanguaging
*Edited by Cenoz Jasone and Gorter Durk*

**Learning Vocabulary in Another Language**
2nd edition
*I. S. P. Nation*

**Narrative Research in Applied Linguistics**
*Edited by Barkhuizen Gary*

**Teacher Research in Language Teaching**
A Critical Analysis
*Simon Borg*

**Figurative Language, Genre and Register**
*Alice Deignan, Jeannette Littlemore and
Elena Semino*

**Exploring ELF**
Academic English Shaped by Non-native
Speakers
*Anna Mauranen*

**Genres across the Disciplines**
Student Writing in Higher Education
*Hilary Nesi and Sheena Gardner*

**Disciplinary Identities**
Individuality and Community in Academic
Discourse
*Ken Hyland*

**Replication Research in Applied Linguistics**
*Edited by PorteGraeme*

**The Language of Business Meetings**
*Michael Handford*

**Reading in a Second Language**
Moving from Theory to Practice
*William Grabe*

**Modelling and Assessing Vocabulary
Knowledge**
*Edited by Daller Helmut, Milton James and
Treffers-Daller Jeanine*

**Practice in a Second Language**
Perspectives from Applied Linguistics and
Cognitive Psychology
*Edited by DeKeyser Robert M.*

**Task-Based Language Education**
From Theory to Practice
*Edited by van den Branden Kris*

**Second Language Needs Analysis**
*Edited by Long Michael H.*

**Insights into Second Language Reading**
A Cross-Linguistic Approach
*Keiko Koda*

**Research Genres**
Exploration and Applications*John M. Swales*

**Critical Pedagogies and Language Learning**
*Edited by Norton Bonny and Toohey Kelleen*

**Exploring the Dynamics of Second Language
Writing**
*Edited by Kroll Barbara*

**Understanding Expertise in Teaching**
Case Studies of Second Language Teachers
*Amy B. M. Tsui*

**Criterion-Referenced Language Testing**
*James Dean Brown and Thom Hudson*

**Corpora in Applied Linguistics**
*Susan Hunston*

**Pragmatics in Language Teaching**
*Edited by Rose Kenneth R. and Kasper Gabriele*

**Cognition and Second Language Instruction**
*Edited by Robinson Peter*

**Research Perspectives on English for
Academic Purposes**
*Edited by Flowerdew John and Peacock
Matthew*

**Computer Applications in Second Language
Acquisition**
Foundations for Teaching, Testing and Research
*Carol A. Chapelle*

# Corpora In Applied Linguistics

## *Susan Hunston*
*University of Birmingham*

CAMBRIDGE
UNIVERSITY PRESS

# CAMBRIDGE
## UNIVERSITY PRESS

For Michael Hoey, 1948–2021

# *Contents*

ix

# *Figures*

xii

*List of Figures*     xiii

# *Tables*

xvi    *List of Tables*

## *Acknowledgements*

This book owes a great deal to the influence of colleagues and fellow researchers in corpus linguistics and applied linguistics. These include friends in the Centre for Corpus Research at the University of Birmingham as well as colleagues around the world. Their support, insight and challenge have been invaluable.

In this book I draw on many of the corpus resources, both corpora and software, provided by the SketchEngine team (sketchengine.eu) and by English Corpora run by Professor Mark Davies (english-corpora.org). Other corpora and resources are acknowledged as appropriate through the book. I am grateful for permission to adapt and reproduce data of various kinds from other published sources, as acknowledged in the relevant places. I would like to mention specifically the John Benjamins series 'Studies in Corpus Linguistics' (benjamins.com/catalog/scl), and the Cambridge University Press series 'Elements in Corpus Linguistics' (www.cambridge.org/core/what-we-publish/elements/corpus-linguistics).

In the 20 years since the publication of the first edition of *Corpora in Applied Linguistics* I have had the privilege of working with PhD students, too many to mention by name, whose work has been inspirational and whose enthusiasm has always brightened my day. I would also like to acknowledge the love and continued support of friends and family, especially Ann, Jyl, Jeannette, Paul, Suganthi, Michaela, Linda and John.