Data and Methods in Corpus Linguistics

Corpus linguistics continues to be a vibrant methodology applied across highly diverse fields of research in the language sciences. With the current steep rise in corpus sizes, computational power, statistical literacy and multi-purpose software tools, and inspired by neighbouring disciplines, approaches have diversified to an extent that calls for an intensification of the accompanying critical debate. Bringing together a team of leading experts, this book follows a unique design, comparing advanced methods and approaches current in corpus linguistics, to stimulate reflective evaluation and discussion. Each chapter explores the strengths and weaknesses of different datasets and techniques, presenting a case study and allowing readers to gauge methodological options in practice. Contributions also provide suggestions for further reading, and data and analysis scripts are included in an online appendix. This is an important and timely volume, and will be essential reading for any linguist interested in corpus-linguistic approaches to variation and change.

OLE SCHÜTZLER is Professor for Varieties of English at Leipzig University. Mostly working within the frameworks of quantitative sociolinguistics/socio-phonetics and corpus linguistics, he takes a general interest in synchronic and diachronic variation and change in English with a special focus on Scottish Englishes.

JULIA SCHLÜTER is Associate Professor for English Linguistics at the University of Bamberg. Her research interests lie in the areas of phonological and grammatical variation in British and American English past and present, empirical – especially corpus-based – methodologies, and applications of linguistic insights and techniques to the teaching of English.

# Data and Methods in Corpus Linguistics

*Comparative Approaches*

*Edited by*

Ole Schützler
*Leipzig University*

Julia Schlüter
*University of Bamberg*

CAMBRIDGE
UNIVERSITY PRESS

CAMBRIDGE
UNIVERSITY PRESS

# Contents

vi      Contents

# Figures

## Tables

---

List of Tables    xiii

xiv    List of Tables

# Contributors

SABINE ARNDT-LAPPE Department of English Studies, University of Trier, Germany

TOBIAS BERNAISCH Department of English, Justus Liebig University Giessen, Germany

DOUGLAS BIBER Department of English, Northern Arizona University, Flagstaff, USA

JESSE EGBERT Department of English, Northern Arizona University, Flagstaff, USA

MATTHEW FAHY Department of Mathematics and Statistics, Northern Arizona University, Flagstaff, USA

VOLKER GAST Department of English and American Studies, University of Jena, Germany

SEBASTIAN HOFFMANN Department of English Studies, University of Trier, Germany

MANFRED KRUG Department of English and American Studies, University of Bamberg, Germany

NATALIA LEVSHINA Neurobiology of Language Department, Max Planck Institute for Psycholinguistics, Nijmegen, The Netherlands

SETH MEHL Digital Humanities Institute, University of Sheffield, United Kingdom

JULIA SCHLÜTER Department of English and American Studies, University of Bamberg, Germany

GEROLD SCHNEIDER Department of Computational Linguistics, University of Zurich, Switzerland

OLE SCHÜTZLER Institute of British Studies, Leipzig University, Germany

LUKAS SÖNNING Department of English and American Studies, University of Bamberg, Germany

BENEDIKT SZMRECSANYI Department of Linguistics, University of Leuven, Belgium

FABIAN VETTER Department of English and American Studies, University of Bamberg, Germany

SEAN WALLIS Survey of English Usage, University College London, United Kingdom

# Acknowledgements