

DE GRUYTER  
MOUTON

*Geoffrey Sampson,  
Anna Babarczy*

# GRAMMAR WITHOUT GRAMMATICALITY

GROWTH AND LIMITS  
OF GRAMMATICAL PRECISION

TRENDS IN LINGUISTICS

Geoffrey Sampson  
Anna Babarczy

# Grammar Without Grammaticality

---

Growth and Limits of Grammatical Precision

**DE GRUYTER**  
MOUTON

ISBN 978-3-11-028977-0  
e-ISBN 978-3-11-029001-1  
ISSN 1861-4302

**Library of Congress Cataloging-in-Publication Data**

A CIP catalog record for this book has been applied for at the Library of Congress.

**Bibliographic information published by the Deutsche Nationalbibliothek**

The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie;  
detailed bibliographic data are available in the Internet at <http://dnb.dnb.de>.

© 2014 Walter de Gruyter GmbH, Berlin/Boston

Typesetting: RoyalStandard, Hong Kong

Printing and binding: Hubert & Co. GmbH & Co. KG, Göttingen

♻️ Printed on acid-free paper

Printed in Germany

[www.degruyter.com](http://www.degruyter.com)



Geoffrey Sampson and Anna Babarczy —  
**Grammar Without Grammaticality**

# **Trends in Linguistics Studies and Monographs**

---

**Editor**

Volker Gast

**Editorial Board**

Walter Bisang

Jan Terje Faarlund

Hans Henrich Hock

Natalia Levshina

Heiko Narrog

Matthias Schlesewsky

Amir Zeldes

Niina Ning Zhang

**Editor responsible for this volume**

Volker Gast

## **Volume 254**

## Preface

This book sets out to explain and give evidence for a distinctive view of the nature of human language, a view which contrasts sharply with assumptions that are taken for granted by many or most linguists publishing today. The book has had a long gestation, but the bulk of it emerges from collaborative work carried out at Sussex University in England, where one of us (Sampson) was Professor of Natural Language Engineering, and the other (Babarczy) was assisting him as Research Fellow.

Some chapters of the book are primarily conceptual in nature, aiming to make as clear and explicit as possible what we think language is like, where we disagree with others, and how we answer their objections to our point of view. Other chapters go into details of the empirical, quantitative researches that have led us to those conceptual conclusions. (Some chapters mix these modes of discussion.)

We realize, of course, that many readers are more comfortable with conceptual discussion than with the nitty-gritty of numbers and statistics. Some might have preferred us to leave out the second kind of material, feeling that they could trust us to get that part right and give them just the conclusions which emerge. We believe that it is important not to omit the quantitative material. There is a style of linguistic publication which is quite influential nowadays, that might be described as a pastiche of scientific writing. General statements are expressed in words, but here and there one sees a sprinkling of algebraic notations or other formalisms which create an impression of hard scientific analysis underlying the generalities – only, if you press, that analysis never seems to be forthcoming. (The kind of intellectual charlatanry skewered by Alan Sokal and Jean Bricmont (1997) is familiar in modern linguistics, although it happened that none of the individual academics Sokal and Bricmont discussed was a linguist.) We want it to be clear to readers that our work is not like that. If some passages of the book seem hard going, we hope that readers will be patient with them for the sake of receiving a sense of the detailed research which, in our eyes at least, has entitled us to assert the conclusions we have drawn.

If some parts of the book may seem dry, in another respect the book may offer *more* human interest than the average work on linguistics. Since language is an aspect of human behaviour, and an important one, any fresh insights on the nature and structure of language will be of some interest to human beings – but often those insights are rather abstract. Linguistics is frequently seen as a subject which does not have the weighty practical or moral implications

for society that commonly attach to fundamental theorizing in other human sciences, such as economics or politics. But we believe that the contrasting views about language structure discussed in this book – ours, and the current majority view with which we disagree – do indeed entail contrasting consequences for how societies should be organized and lives should be led. The central aim of our book is to establish a truth about language, but professional academics have a duty to society to spend some time discussing the practical implications of their work. We take this issue up in our final chapter.

Some of our chapters are based on papers previously published in academic journals or conference proceedings, with redrafting and additional material to fit them into a volume that reads coherently as a whole. Only in this way can separate small-scale investigations be appreciated as contributing to a single novel and consistent model of the nature of language. Other chapters are newly published here.

Where material was published previously, this was sometimes done under one of our names, sometimes the other, and sometimes we published as co-authors. That creates a problem about pronouns for this book. It would be tedious for readers if we chopped and changed between “I” and “we”, depending on details of the original published version of different passages here. Furthermore, even when Sampson published under his sole name, often his writing drew on work done by Babarczy.

Luckily, in English the authorial “we” is vague between singular and plural: for the sake of readability, we shall be “we” throughout this book, even in contexts where the reference is obviously singular. For instance, at the beginning of chapter 1 we mention “our copy” of an old book by John Meiklejohn: of course the book belongs to only one of us, and we have separate bookshelves, now in different countries – but we spare the reader the need to wade through irrelevant references to “one of us”, “the other co-author”, and so forth.

The exception is that at a few points it matters which of us is which. For instance, in chapter 2 we describe an experiment in which we played separate roles and where it could have been relevant that GRS but not AB is an English native speaker. In such contexts we identify ourselves by initials, as we have just done. Broadly, as will commonly be the case when a senior academic collaborates with a full-time researcher, the bulk of the hard quantitative work tended to be done by AB, and the general conclusions, together with the ways to make them engage with the ongoing context of debate that constitutes the discipline of linguistics, tended to be settled by GRS. But each of us contributed at each level; the book is a genuine collaboration.

The programme of research from which this book emerges has centred round various English-language “corpora” – samples, held in electronic form,

of different genres of the language. Since many chapters will refer to one or another of these resources, we list here those which are relevant to our book, with access details. (Sampson and McCarthy 2004 contains papers discussing these and many other corpora, and illustrating various ways in which they are used for linguistic research.)

The “raw” English corpora mainly relevant to our work – where “raw” means that the electronic files contain the wording of the language samples but little more<sup>1</sup> – are the *Brown Corpus*, the *Lancaster–Oslo/Bergen (LOB) Corpus*, and the *British National Corpus*. Brown and LOB are the resources which first enabled corpus linguistics to get off the ground: they are matched collections of one million words each of published writing in, respectively, American and British English of the 1960s. The British National Corpus is a more recent collection comprising one hundred million words of British English, of which ninety million is written language but ten million is spoken, and of the latter about 4.2 million words are “demographically sampled” speech – broadly, the spontaneous conversational talk of a cross-section of the UK population. For Brown and LOB Corpora see e.g. <khnt.hit.uib.no/icame/manuals/>; for the British National Corpus, <www.natcorp.ox.ac.uk>.

Initially as a member of the corpus research group at Lancaster University led by Geoffrey Leech, and later as leader of his own research group of which AB became a key member, GRS has used these and similar materials in order to generate a maximally refined and comprehensive structural taxonomy for written and spoken English, attempting to do for the grammar of the English language something akin to what Linnaeus in the eighteenth century did for plant life. (On this analogy, cf. Sampson 2000a.) In practice this entailed using subsamples of the raw corpora as testbeds to uncover and resolve shortcomings in a scheme of structural annotation for the language, making the scheme ever more adequate to record and classify the grammatical detail that occurs in samples of the language in practice – a process which can never be complete. The resulting annotated subcorpora, or “treebanks”, then became research resources in their own right:

- the SUSANNE Corpus comprises 64 files drawn from four genre categories of the Brown Corpus, for a total of about 130,000 words of written American English
- the CHRISTINE Corpus comprises about 80,500 words, plus many ums and ers and the like, of spoken British English (taken from the demographically-sampled speech section of the British National Corpus)

---

<sup>1</sup> In fact, for these resources, recent versions of the files do also include wordtags, i.e. part-of-speech classifications of the individual words.



- the LUCY Corpus comprises 165,000 words of written British English representing a spectrum of genres ranging from polished published writing to the writing of young children.

(The name SUSANNE was chosen to stand for “surface and underlying structural analyses of natural English”, and also to celebrate links with the life of St Susanna. The names CHRISTINE and LUCY likewise refer to saints whose reputed careers were in different ways appropriate to the language genres in question.)

These research resources are freely available for others to download and work with (and many researchers internationally have done and are doing so). For details of access, and documentation of these resources, see <[www.grsampson.net/Resources.html](http://www.grsampson.net/Resources.html)>.

The structural annotation scheme which emerged from this programme of research, called the SUSANNE Scheme after the earliest of the treebanks listed, was published in the form of a 500-page book (Sampson 1995). Various refinements to that published scheme which derived from work since 1995, particularly to handle the special features of spontaneous speech, are covered in the electronic documentation files for the corpora.

Geoffrey Sampson  
University of South Africa

Anna Babarczy  
Budapest University of Technology and Economics

# Acknowledgements

We thank Alan Morris for his work on the electronic resources used in the research discussed in this book, and Anna Rahman for permission to use the material on which chapter 6 is based, of which she was a co-author. For discussion of various of our ideas, we are grateful to Anders Ahlqvist, John Carroll, Gerald Gazdar, Adam Kilgariff, and the late Larry Trask, as well as numerous participants at conferences and anonymous referees for journals. We ask anyone whose name we have overlooked to forgive us.

Much of the research presented here was sponsored by the Economic and Social Research Council (UK), whose support we gratefully acknowledge.

The origin of the various chapters is as follows:

Chapter 1 is new material.

Chapter 2 is based on a paper in the *Journal of Natural Language Engineering*, vol. 14, pp. 471–494, 2008, reprinted with permission of Cambridge University Press.

Chapter 3 is based on a paper in the *Proceedings of the Workshop on Linguistically Interpreted Corpora, LINC-2000*, Luxemburg, 6 Aug 2000 (ed. by Anne Abeillé et al.), pp. 28–34.

Chapter 4 is based on a “target paper” which appeared in *Corpus Linguistics and Linguistic Theory*, vol. 3, pp. 1–32 and 111–129, 2007.

Chapter 5 is new material.

Chapter 6 is based on a paper in J. M. Kirk (ed.), *Corpora Galore*, 2000, pp. 295–311, reprinted by permission of Rodopi of Amsterdam.

Chapter 7 is based on a chapter in *Empirical Linguistics*, © Geoffrey Sampson 2001, reprinted by permission of Continuum Publishing, an imprint of Bloomsbury Publishing plc.

Chapter 8 is based on a paper in Sylviane Granger and Stephanie Petch-Tyson (eds), *Extending the Scope of Corpus-Based Research*, 2003, pp. 177–193, reprinted by permission of Rodopi of Amsterdam.

Chapter 9 is new material.

Chapter 10 is based on a paper in Mickael Suominen et al. (eds), *A Man of Measure: Festschrift in honour of Fred Karlsson on his 60th birthday*, 2006, pp. 362–374, reprinted with permission of the Linguistic Association of Finland.

Chapter 11 is based on a paper in *English Language and Linguistics*, vol. 6, pp. 17–30, 2002, reprinted with permission of Cambridge University Press.

Chapter 12 is based on a paper in the *Journal of Natural Language Engineering*, vol. 9, pp. 365–380, 2003, reprinted with permission of Cambridge University Press.

Chapter 13 is based on a paper in the *International Journal of Corpus Linguistics*, vol. 10, pp. 15–36, 2005.

Chapter 14 is new material.

Chapter 15 is based on a keynote address to the International Association for Dialogue Analysis workshop (IADA 2006), Mainz, September 2006, published in Marion Grein and Edda Weigand (eds), *Dialogue and Culture*, pp. 3–25, 2007, and reprinted with kind permission of John Benjamins Publishing Company.

## List of figures

- Figure 1: Meiklejohn's structural diagram of a William Morris couplet — 3
- Figure 2: Meiklejohn's diagram recast in tree format — 3
- Figure 3: A SUSANNE tree structure — 32
- Figure 4: Locations of analytic discrepancy — 39
- Figure 5: Causes of analytic discrepancy — 42
- Figure 6: Noun-phrase expansion frequencies in SUSANNE — 71
- Figure 7: Tree structure of a CHRISTINE speaker turn — 121
- Figure 8: Embedding scores of words in a CHRISTINE extract — 140
- Figure 9: Embedding indices plotted against age — 150
- Figure 10: Frequencies of clause types in conducive contexts, by genre — 176
- Figure 11: Frequencies of clause types by clause-function of anaphor and by genre — 177
- Figure 12: Frequencies of clause types by phrase-function of anaphor and by genre — 179
- Figure 13: Frequencies of clause types by antecedent position and by genre — 180
- Figure 14: GEIG/unlabelled v. LA parse scores — 226
- Figure 15: GEIG/labelled v. LA parse scores — 227
- Figure 16: Corpus-based computational linguistics papers (after Hirschberg 1998) — 254
- Figure 17: Raw 3-way classification of *Language* articles — 261
- Figure 18: Evidence-based as % of non-neutral *Language* articles (smoothed) — 263

## List of tables

- Table 1: Complexity figures for the sample texts — 35
- Table 2: Inter-analyst agreement on sample text analyses — 36
- Table 3: Results of discrepancy resolution — 47
- Table 4: Mean embedding indices by region — 144
- Table 5: Mean embedding indices by social class — 145
- Table 6: Mean embedding indices by sex — 146
- Table 7: Mean embedding indices by age — 146
- Table 8: Mean adult embedding indices by sex — 147
- Table 9: Mean adult embedding indices by age — 147
- Table 10: Incidence of phrase categories in different genres — 164
- Table 11: Incidence of subordinate-clause categories in different genres — 166
- Table 12: Incidence of relative-clause types in different genres — 171
- Table 13: Effects of context type on relative-clause selection — 182
- Table 14: Verb qualifier frequencies — 206
- Table 15: Perfect and Past marking by region — 207
- Table 16: *got* for standard *HAVE got* by region — 209
- Table 17: Ranking on either metric of parses ranked lowest by the other metric — 230
- Table 18: Leaf-ancestor scores for individual words in a parse tree — 235

# Table of contents

Preface — v

Acknowledgements — ix

List of figures — xvi

List of tables — xvii

## 1 Introduction — 1

- 1.1 Grammar before linguistics — 1
- 1.2 All grammars leak — 6
- 1.3 No common logic — 8
- 1.4 “Chicken eat” — 9
- 1.5 The case of Old Chinese — 13
- 1.6 It cuts both ways — 18
- 1.7 Vocabulary differences — 19
- 1.8 What can be said about grammar — 21
- 1.9 The computational viewpoint — 24

## 2 The bounds of grammatical refinement — 26

- 2.1 An experiment — 26
- 2.2 The experimental material — 29
- 2.3 The analytic scheme — 31
- 2.4 Measuring similarity of analyses — 34
- 2.5 Text complexity — 35
- 2.6 Overall similarity results — 36
- 2.7 Dividing overall discrepancy between annotation categories — 37
- 2.8 Assigning responsibility for discrepancies — 40
- 2.9 Monitoring for bias — 46
- 2.10 Implications of the experiment — 47
- 2.11 New research techniques yield novel perspectives — 52

## 3 Where should annotation stop? — 53

- 3.1 Another way to survey indeterminacy — 53
- 3.2 Detailed v. skeleton analytic schemes — 53
- 3.3 The trainability criterion — 55
- 3.4 Limits to expert decision-making — 55
- 3.5 Some examples of indeterminacy — 56
- 3.6 Annotation practice and linguistic theory — 61
- 3.7 A disanalogy with biology — 62

<b>4</b>	<b>Grammar without grammaticality — 64</b>
4.1	Strangers or unmet friends — 64
4.2	Unfamiliar does not imply ungrammatical — 66
4.3	Statistics of construction frequencies — 70
4.4	A range without boundaries — 75
4.5	Can intuition substitute for observation? — 78
4.6	How intuitions have misled — 81
4.7	Is English special? — 84
4.8	The analogy with word meaning — 86
4.9	Grammar as an expression of logical structure — 90
4.10	Realistic grammatical description — 92
<b>5</b>	<b>Replies to our critics — 95</b>
5.1	Is our idea controversial? — 95
5.2	Geoffrey Pullum's objections — 96
5.3	No virtue in extremism — 98
5.4	Stefanowitsch versus Müller — 99
5.5	Trees have no legs — 101
5.6	Law versus good behaviour — 103
5.7	Conceptual objections to our thesis — 104
5.8	Do we really mean it? — 105
5.9	Grammaticality implied by Universal Grammar — 107
5.10	The downfall of Universal Grammar — 109
5.11	Economic growth and linguistic theory — 112
5.12	Discipline should not contradict discipline — 116
5.13	Language is not "special" — 117
<b>6</b>	<b>Grammatical description meets spontaneous speech — 119</b>
6.1	The primacy of speech — 119
6.2	An example — 120
6.3	Wordtagging — 123
6.4	Speech repairs — 124
6.5	Syntactically Markovian constructions — 125
6.6	Logical distinctions dependent on the written medium — 127
6.7	Nonstandard usage — 128
6.8	Dialect difference versus performance error — 130
6.9	Transcription inadequacies — 132
6.10	Dropping the paradigm — 133

<b>7</b>	<b>Demographic correlates of speech complexity — 136</b>
7.1	Speech in the British National Corpus — 136
7.2	Measuring speech complexity — 138
7.3	Classifying the speakers — 141
7.4	Demographics and complexity indices compared — 144
7.5	“Critical period” or lifelong learning? — 148
7.6	Individual advance or collective retreat? — 153
<b>8</b>	<b>The structure of children’s writing — 155</b>
8.1	Moving from spoken to adult written norms — 155
8.2	The language samples — 155
8.3	The suitability of the child-writing sample — 156
8.4	Writing “wordier” than speech — 157
8.5	Width v. depth in parse-trees — 158
8.6	Interim summary — 161
8.7	Phrase and clause categories — 162
8.8	Use of phrase categories — 163
8.9	Use of subordinate clause categories — 165
8.10	The complexity of the relative constructions — 168
8.11	Simple v. complex relatives — 169
8.12	Unanswered questions — 171
<b>9</b>	<b>Child writing and discourse organization — 172</b>
9.1	A fixed grammatical programme? — 172
9.2	New information about a previously identified object — 172
9.3	The new study: data and methods of analysis — 173
9.4	Context frequency — 175
9.5	Syntactic patterns — 176
9.6	Mistakes with relative clauses — 179
9.7	The upshot of the analysis — 182
<b>10</b>	<b>Simple grammars and new grammars — 184</b>
10.1	Pidgins and creoles — 184
10.2	Old Chinese as a counterexample — 187
10.3	Old Chinese not a creole — 187
10.4	Examples of structural vagueness — 188
10.5	Lack of word classes — 190
10.6	Logical indeterminacy — 191
10.7	McWhorter’s diagnostics — 193
10.8	No tone in Old Chinese — 193
10.9	No inflexion in Old Chinese — 194



10.10	Derivational morphology in Old Chinese —	194
10.11	An accident of history —	196
10.12	“Hidden” versus “overt” structure —	197
10.13	Deutscher on Akkadian —	199
10.14	Diverse paths of evolution —	199
<b>11</b>	<b>The case of the vanishing perfect —</b>	<b>201</b>
11.1	Losses as well as gains —	201
11.2	The Perfect aspect and spontaneous speech —	202
11.3	The standard system and nonstandard alternatives —	203
11.4	Verb qualifiers in CHRISTINE —	205
11.5	Past and Perfect —	207
11.6	<i>got</i> for <i>HAVE got</i> —	209
11.7	Casual subject-auxiliary omission —	209
11.8	Modals + <i>of</i> —	210
11.9	Nonstandard verb forms —	212
11.10	A possible explanation —	214
11.11	If one feature can go, what cannot? —	217
<b>12</b>	<b>Testing a metric for parse accuracy —</b>	<b>218</b>
12.1	The need for a metric —	218
12.2	Alternative metrics —	219
12.3	The essence of leaf-ancestor assessment —	220
12.4	The experimental material —	222
12.5	Calculation of lineage similarity —	224
12.6	Are the metrics equivalent? —	226
12.7	Performance systematically compared —	229
12.8	Local error information —	234
12.9	Authority is fallible —	236
<b>13</b>	<b>Linguistics empirical and unempirical —</b>	<b>237</b>
13.1	What went wrong? —	237
13.2	Two kinds of empiricism —	237
13.3	Universal Grammar versus empiricism —	239
13.4	Arguments against empiricism —	241
13.5	How empirical should linguistics be? —	243
13.6	How intuition has led linguists astray —	243
13.7	Were our intuitions correct after all? —	246
13.8	Can intuitions be empirical? —	248
13.9	Is our characterization of generative linguistics misleading? —	251