



The Encyclopedia of Applied Linguistics

Corpora in World Englishes (General issues in World Englishes)

Journal:	<i>The Encyclopedia of Applied Linguistics</i>
Manuscript ID	wbeal20127.R1
Wiley - Manuscript type:	Encyclopedia of Applied Linguistics article
Date Submitted by the Author:	n/a
Complete List of Authors:	Gries, Stefan Th.; University of California Santa Barbara, Linguistics; Justus-Liebig-Universitat Giessen, English Deshors, Sandra; no affiliation
Keywords:	corpora, varieties research, nativization, epicenter research, statistical modeling

SCHOLARONE™
Manuscripts

Corpora in World Englishes

Sandra C. Deshors
no affiliation

Stefan Th. Gries
UC Santa Barbara & JLU Giessen
stgries@linguistics.ucsb.edu

Abstract

In this overview, we survey recent and current developments in corpus-based research on World Englishes. We exemplify current strands of research in both more theoretical and more applied parts of research on varieties of English and conclude with theoretical, methodological, and resource desiderata.

Key words

corpora, varieties research, nativization, epicenter research, statistical modeling

Introduction

Over the past fifty years the field of World Englishes (WEs) has undergone substantial methodological and theoretical developments that have gone hand-in-hand. A particularly influential methodological development has been the widespread turn toward corpora (singular, *corpus*), i.e., collections of naturally occurring spoken or written text in electronic format. This adoption of corpus methods has been accompanied by a growing adoption of modern statistical approaches, which have informed studies of concrete phenomena (e.g., grammatical variation across Englishes) or the development of theoretical models of world Englishes. As Deshors & Bernaisch (2019: 85) note, through methodological advancements “not only have WE researchers managed to better understand the linguistic forces that drive the development of Englishes worldwide, but corpus-based research of world Englishes has become its own academic sub-field”. Two very broadly-defined kinds of studies can be distinguished: those that are more theoretical and structural in nature, and those whose focus is more applied and anchored in sociolinguistic and cultural frameworks. In what follows, we discuss these two trends in that order and provide methodological commentary. In the final part of our chapter, we briefly discuss current developments in corpus-based WEs research and identify main desiderata for future work.

More theoretical kinds of research

WEs Scholars with more theoretical research agendas have used corpora to explore a wide range of linguistic areas including language structure (especially, phonology, morpho-syntax, lexicogrammar), but also semantics and cross-variety linguistic variation. Much such work is anchored in usage-based frameworks (e.g., Construction Grammar), which assume that speakers’ knowledge of linguistic items is correlated with the items’ distributional characteristics in

1
2
3 authentic language (e.g., Langacker, 1987; Goldberg, 2006). Rautioaho, Deshors, & Meriläinen
4 (2018) apply this perspective to the progressive vs. non-progressive alternation by exploring
5 three main types of Englishes from Kachru's (1989) classification of Englishes:
6
7

- 8 • English as second language (ESL): Indian, Singaporean and Nigerian Englishes from the
9 *International Corpus of English (ICE)*;
- 10 • English as foreign language (EFL): Finnish, French, and Polish learner Englishes from
11 the *Louvain International Database of Spoken English Interlanguage* and the
12 *International Corpus of Learner English*;
- 13 • native Englishes (ENL): British and American Englishes from ICE and the *Santa*
14 *Barbara Corpus of Spoken American English*.

15
16
17 To explore how the alternation is conditioned by different contexts across English types,
18 the authors use state-of-the-art statistical approaches (cluster analysis and logistic regression
19 modeling) and find that individual varieties conform well to the traditional Kachruvian model:
20 With non-progressives especially, there is a clear-cut divide between ENLs, ESLs, and EFLs.
21 However, they caution that the validity of the ENL-ESL-EFL continuum can be influenced by
22 scholars' choice of quantitative approach, thereby underscoring the tight connection between
23 theory and method in WEs' corpus-based research.
24

25 Also exploring the classification of Englishes via morpho-syntax, Szmrecsanyi &
26 Kortmann (2011) sets high methodological standards for the field. Using cluster-analytic
27 techniques, they demonstrate the impact that sophisticated quantitative corpus-based methods
28 can have on our understanding of the typology of Englishes world-wide and how structural (dis-
29)similarities across Englishes are best explained based on the combined effects of geography,
30 type of English (i.e., EFL, ESL, ENL) and linguistic features. This type of approach helps
31 identify typological issues/features (e.g., lower/higher degrees of morpho-syntactic or structural
32 complexity) across foreign- and second-language varieties, which may have important, applied
33 pedagogical implications for second language learners, instructors, and users.
34

35 Schützler (2020) is a good example of how sophisticated statistical techniques (here,
36 Bayesian mixed-effect logistic regression modeling) to corpora of world Englishes can lead to
37 insightful discoveries on inter-varietal differentiation in the area of lexicogrammar. He explores
38 the variation of the clausal positioning of 1,259 *although*-constructions (final vs. non-final) based
39 on language production mode (speech vs. writing), cognitive processing and semantics across
40 British, Canadian, New Zealand, Nigerian, Indian, and Philippine English (from the ICE). Across
41 all varieties, *although*-constructions prefer final positioning in speech and in dialogic contexts
42 and "cognitive, processing and production-based constraints are strong enough to keep inter-
43 varietal differentiation in check" (Schützler 2020: 455).
44

45
46 Corpus-based theoretical studies in phonology are relatively rare. However, Hay et al.
47 (2018) shows the potential of corpora to unveil, understand, and theorize the phonological
48 patterns that characterize and distinguish Englishes world-wide. In two small corpus-based case
49 studies they explore the usage patterns of linking- and intrusive-*r* usage in (nonrhotic) New
50 Zealand English (NZE) to assess whether the patterns characteristic of early NZE are still
51 observed today, based on both spontaneous speech and speech from a reading passage. The
52 mixed-effects logistic regression statistical approach identifies (i) developmental patterns of
53 linking- and intrusive-*r* through time, (ii) influential patterns of morpheme and word boundaries
54 on the production of *r*-sandhi and (iii) different usage patterns of *r*-sandhi by men and women.
55
56
57

Corpora have also allowed scholars to explore semantic questions in innovative ways. Mehl (2019) is a study on variation focusing on the three light verbs *make*, *take* and *give* and light verb constructions (LVCs) in Singapore, Hong Kong and British English (from the ICE). What stands out in the approach is the unusual direction in which form and meaning are studied, namely the onomasiological perspective from-meaning-to-form as opposed to from-form-to-meaning. Mehl (2019) finds that semantic patterns are constant across Englishes and notes, “there is no evidence for unique or innovative LVCs in the three corpora [and] there is remarkable similarity across the three corpora as well: all varieties, in most cases, prefer the related verb over the LVC in both speech and writing” (Mehl 2019: 77).

Finally, in the area of language variation, corpora have helped identify ways in which individual dialects of native English can, overtime, influence the development of a non-native English variety. For example, Borlongan (2021) is a diachronic (1960s vs. 2000s) study of the understudied variety of Japanese English (in the Diachronic Corpus of Expanding Circle Englishes) and its use of putative American variants of spellings, single words, compound lexical items, lexical endings, suffixes in compounds, verb morphology, contractions, article usage, constituent sequence, verb-noun collocations, preposition choices, collective noun concord and phraseology. Overall, the study portrays Japanese English as “overwhelmingly American in its choices of variants across the categories under investigation” (Borlongan 2021: 54).

In sum, while the above overview is necessarily brief, it is clear that corpus-linguistic analysis has an increasingly greater impact on just about all areas of WEs research.

More applied kinds of research

Sociolinguistic, pragmatic or applied linguistics research has similarly benefited from the rise of corpus methods in linguistics and in WEs research in particular. The interconnectedness of linguistic, social, pragmatic, cultural, multilingual, and communicative aspects of language use is recognized as “a driving force beyond the structural development of Englishes” (Deshors & Gilquin 2018: 287). A first case in point is van Rooy & Kruger (2018). Although that study is primarily a theoretical discussion on how to expand on current theoretical models of WEs, it also shows how multilingual digital repertoires can help us address empirical challenges resulting from the globalization of Englishes such as multilingualism, hybrid varieties, online communication, and complex identities. Based on a corpus of interactive user comments online that accompany daily summaries of the content of the most popular South African television soap operas, Rooy & Kruger (2018) show how the core of online interactions consists of a shared pool of English (lexical) resources and global as well as local nonstandard English forms complemented by forms from South African languages. Ultimately, Rooy & Kruger (2018) showcase the importance of accounting for world Englishes sociolinguistically and in ways that are ecologically valid for developing theoretical models of Englishes world-wide.

Schneider (2018), Funke & Bernaisch (2022), and Revis & Bernaisch (2020) are also located at the interface of culture and corpora and allow us to better understand who the speakers of WEs are and how their identities stand out systematically in corpus data. Starting with Schneider (2018), the study asks, how and to what extent traces of cultural impact can be detected in corpora. The study uses five components of the ICE representing (i) English as spoken in these countries and (ii) major cultural traditions (specifically, Great Britain for a western culture; Hong Kong, Singapore, and India for Asian cultures; Nigeria for West African

1
2
3 cultures). ICE data were occasionally supplemented with data from the Corpus of Global Web-
4 Based English and quantitatively explored along three aspects: linguistic forms referring to
5 objects/artefacts, expressions of cultural dimensions (e.g., collectivism vs. individualism, power
6 distance, social relations, etc.), and linguistic constructions. Maybe unsurprisingly, it emerged
7 that terms for concrete objects, artefacts, and notions appear in regional corpora (mainly locally)
8 but cultural dimensions, while recognizable in corpus data, “vary from one dimension to another
9 and by indicator terms, and they tend to be graded rather than absolute” (Schneider 2018: 128).
10 However, with linguistic constructions, traces of cultural influence are less clear: “such
11 influences exist and have an impact, but if so, they are clearly more indirect and somewhat
12 abstract” (Schneider 2018: 128). Overall, the study shows how pervasive and diverse cultural
13 information can be in corpora WEs and how informative such data are for WEs research.
14
15

16 Finally, and as for discourse and pragmatics, Revis & Bernaisch (2020) investigate how
17 corpora inform discussions on pragmatic nativisation across Englishes, a topic that remains today
18 rarely explored. They focus on filled (e.g., *uh* or *uhm*) and unfilled pauses (i.e., silence) by
19 speakers of Indian and Sri Lankan English (in the ICE) and explore whether (i) variety-specific
20 differences exist in the use of filled/unfilled pauses and (ii) what factors influence speakers’
21 choice of pauses – specifically, whether speakers are influenced by the internal structure of texts
22 and whether their choices are speaker-related, genre-related or related to speakers’ English
23 variety. They include a wide variety of structural and socio-biographical factors to their analysis
24 (e.g., word class, dialogue vs. monologue, scripted vs. unscripted text, speaker’s age, gender, and
25 English variety) and, therefore, use advanced predictive modeling techniques (conditional
26 inference trees and generalized linear mixed-effects models). Their results show that (i) pauses
27 cluster, (ii) their choices are influenced by their frequencies in texts and (iii) filled pauses occur
28 less in monologues. Funke & Bernaisch (2022) is a methodologically similar study – also
29 adopting a multifactorial predictive modeling methods (this time, random forest) – to explore
30 how socio-biographic and pragmatic factors contribute to the use of intensifiers and downtoners
31 also in Indian and Sri Lankan Englishes. Broadly, across varieties, intensifiers and downtoners
32 are more similar than different and they are used more frequently by younger female in informal
33 conversations.
34
35

36 As before, our discussion had to be selective, but hopefully still illustrates the wide
37 variety of corpus methods and statistical applications in this part of WEs research.
38
39
40

41 **Desiderata for the future**

42
43 By definition, corpus-linguistic research on WEs is where varieties research and general corpus
44 linguistics intersect, which means its desiderata are motivated from both these disciplines. As for
45 the former, the discipline of WEs needs to reconcile how multiple forces are pulling it into
46 different directions. On the one hand, the field needs to address what its models are trying to do,
47 what counts as important phenomena, and what counts as evidence. For instance, how much
48 predictive capability do we require our models to have and how much do we allow them to be
49 falsified from different kinds of data? (See Bernaisch et al. 2022 for discussion and critique of
50 several common assumptions.) Relatedly, do we focus on, to use their terminology, linguistic
51 butterflies (often infrequent surface structure deviations from, e.g., British English) or linguistic
52 ants (often frequent structural choices with higher functional loads)? And, finally, how do we
53 integrate usually current linguistic data with often diachronic sociocultural information (e.g.,
54
55
56
57
58
59
60

attitudinal and/or sociolinguistic information) and how do we reconcile different explanations for the same empirical findings?

On the other hand, there are many open methodological questions, some of them involving statistical, others involving resource availability. For example, how do we study what we study – do we rely on observed (relative) frequencies across varieties including comparisons with the presumed source variety or do we rely (more) on multifactorial statistical modeling and exploratory tools? And if we do the latter, what kinds of predictors are most relevant? Not all of the field has realized that even just for statistical reasons alone, level-1 predictors (i.e., observation-level predictors describing each individual speaker choice) need to be included to be able to for any generalizations regarding level-2 predictors (e.g., speaker-specific predictors) or level-3 predictors (e.g., variety differences) to be valid (see Gries, 2023, Section 2). Lastly, we need better resources; most importantly, we need

- more corpora with more diverse register/genre coverage to better study the dispersion of features – ants and butterflies – in and across varieties;
- more diachronic corpora for research on phenomena such as epicentral influence, because Gries et al. (2018) have demonstrated that the currently predominant apparent-time approach to indigenization or nativization phenomena cannot deliver the results it has been claiming to deliver.

With such additions to our methodological toolkit and corpus resources and a renewed focus on what theoretical models can and should do, corpus-based research on WEs will continue to evolve in promising ways.

SEE ALSO: World Englishes and the Native Speaker; The Dynamic Model of Postcolonial English; Edgar W. Schneider; Bilingualism and Multilingualism – an overview of the field; Corpus Linguistics: Overview; Corpus Linguistics: Quantitative Methods

References

- Bernaisch, T. J. Gries, St. Th., & Heller, B. (2022). On the relation between theoretical models and statistical modeling: the case of linguistic epicentres. *World Englishes*, 41, 333-346.
- Borlongan, A. M. (2021). A new American-lineage English? Proportions of American variants in Japanese English}. *Asian Englishes*, 23, (1), 51-61.
- Deshors, S. C. & Bernaisch, T. (2019). Corpus approaches to World Englishes: A bird's eye view. In P. de Costa, D. Crowther, & J. Maloney (eds.), *Investigating World Englishes: Research Methodology and Practical Applications* (pp. 21-43). New York: Routledge.
- Deshors, S. C. & Gilquin, G. (2018). Modeling world Englishes in the 21st century: New reflections on model-making. In S. C. Deshors (ed.), *Modeling World Englishes in the 21st Century: Assessing the Interplay of Emancipation and Globalization of ESL Varieties*, (pp. 281-294). Amsterdam & Philadelphia: John Benjamins.
- Funke, N. & Bernaisch, T. J. (2022). Intensifying and downtoning in South Asian Englishes: Empirical perspectives. *English World-Wide*, 43, (1), 33-65.

- 1
2
3 Goldberg, A. (2006). *Constructions at work: The nature of generalization in language*. Oxford:
4 Oxford University Press.
- 5 Gries, St. Th. (2023). Corpus-linguistic and computational methods for analyzing communicative
6 competence: contributions from usage-based approaches. In M.H. Kanwit & M. Solon
7 (eds.), *Communicative competence in a second language: theory, method, and*
8 *applications* (pp. 115-131). New York & London: Routledge.
- 9 Gries, St. Th., Bernaisch, T. J., & Heller, B. (2018). A corpus-linguistic account of the history of
10 the genitive alternation in Singapore English. In S. C. Deshors (ed.), *Modeling World*
11 *Englishes: Assessing the interplay of emancipation and globalization of ESL varieties*,
12 (pp. 245-279). Amsterdam & Philadelphia: John Benjamins.
- 13 Haselow, A. (2021). Dealing with trouble in conversation in English-speaking cultures:
14 Conversational repair in global varieties of English. *English World-Wide*, 42, (3), 324-
15 349.
- 16 Hay, J., Drager, K., & Gibson, A. (2018). Hearing R-sandhi: The role of past experience.
17 *Language*, 94, (2), 360-404.
- 18 Kachru, B. B. (1989). World Englishes and applied linguistics. In M. L. Tickoo (ed.), *Languages*
19 *& Standards: Issues, Attitudes, Case Studies*, (pp. 178-205). Singapore: SEAMEO
20 Regional Language Centre.
- 21 Kortmann, B. (2010). Variation across Englishes: Syntax. In A. Kirkpatrick (Ed.), *The Routledge*
22 *Handbook of World Englishes* (pp. 422-446). United Kingdom: Routledge.
- 23 Kortmann, B. & Lunkenheimer, K. & Ehret, K. (Eds.). (2020). *The Electronic World Atlas of*
24 *Varieties of English*. Zenodo.
- 25 Kruger, H. & van Rooy, B. (2019). A multifactorial analysis of contact-induced change in speech
26 reporting in written White South African English (WSAfE). *English Language &*
27 *Linguistics*, 24, (1), 179-209.
- 28 Laitinen, M. (2020). Empirical perspectives on English as a lingua franca (ELF) grammar. *World*
29 *Englishes*, 39, (3), 427-442.
- 30 Langacker, R. (1987). *Foundations of cognitive grammar: Theoretical prerequisites* (Vol. 1).
31 Stanford, CA: Stanford University Press.
- 32 Mehl, S. (2019). Light verb semantics in the International Corpus of English: Onomasiological
33 variation, identity evidence and degrees of lightness. *English Language and Linguistics*,
34 23, (1), 55-80.
- 35 Rautionaho, P, Deshors, S. C., & Meriläinen, L. (2018). Revisiting the ENL-ESL-EFL
36 continuum: A multifactorial approach to grammatical aspect in spoken Englishes, *ICAME*
37 *Journal*, 42, 31-68.
- 38 Revis, M. & Bernaisch, T. J. (2020). The pragmatic nativisation of pauses in Asian Englishes.
39 *World Englishes*, 39, (1), 135-153.
- 40 van Rooy, B. & Kruger, H. (2018). Hybridity, globalisation and models of Englishes: English in
41 South African multilingual digital repertoires. In S. C. Deshors (ed.), *Modeling World*
42 *Englishes: Assessing the interplay of emancipation and globalization of ESL varieties*,
43 (pp. 77-108). Amsterdam & Philadelphia: John Benjamins.
- 44 Schneider, E. W. (2018). The interface between cultures and corpora: Tracing reflections and
45 manifestations. *ICAME Journal*, 42, 97-132.
- 46 Schützler, O. (2020). *Although*-constructions in varieties of English. *World Englishes*, 39, (3),
47 443-461.
- 48 Szmrecsanyi, B. & Kortmann, B. (2011). Typological profiling: Learner Englishes vs.

1
2
3 Indigenized L2 varieties of English. In J. Mukherjee & M. Hundt (eds), *Exploring*
4 *Second-Language Varieties of English and Learner Englishes: Bridging a paradigm gap*
5 (pp. 167-188). Amsterdam: John Benjamins.
6
7

8 9 **Suggested Readings**

- 10
11 Gries, St. Th. (2016). *Quantitative corpus linguistics with R: A practical introduction*. 2nd ed.
12 New York: Routledge.
13 Schreier, D., Hundt, M. & Schneider, E. W. (eds.) (2020). *The Cambridge Handbook of World*
14 *Englishes*. Cambridge: Cambridge University Press.
15 Stefanowitsch, A. (2020). *Corpus linguistics: A guide to the methodology*. Berlin: Language
16 Science Press.
17
18
19

20 21 **Contributor Bios**

22
23 Sandra C. Deshors is Associate Professor in the Department of Linguistics, Languages, and
24 Cultures at Michigan State University where she is a core faculty member in the Second
25 Language Studies Ph.D. Program and the Master of Arts in TESOL Program. Her research is
26 anchored in the usage-based theoretical framework and focusses on quantitative corpus-based
27 approaches to learner language, English as a Second Language, and World Englishes.
28

29
30 Stefan Th. Gries is a Professor in the Department of Linguistics at the University of California,
31 Santa Barbara and Chair of English Linguistics (Corpus Linguistics with a focus on quantitative
32 methods) in the Department of English at the Justus-Liebig-University of Giessen. He does
33 research on quantitative corpus linguistics and statistical methods in (corpus) linguistics from a
34 usage-based theoretical perspective.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57

Cross References:

World Englishes and the Native Speaker

The Dynamic Model of Postcolonial English

Edgar W. Schneider

Bilingualism and Multilingualism - an overview of the field

Corpus Linguistics: Overview

Corpus Linguistics: Quantitative Methods

For Peer Review

1
2
3 Stefan Th. Gries is a Professor in the Department of Linguistics at
4 the University of California, Santa Barbara (UCSB) and since 2018 also
5 Chair of English Linguistics (Corpus Linguistics with a focus on
6 quantitative methods, 25%) at the Justus-Liebig-Universität Giessen.
7 Gries earned his Ph.D. degree at the University of Hamburg in 2000. He
8 worked at the University of Southern Denmark at Sønderborg
9 (1998-2005), first as a lecturer, then as assistant professor and
10 tenured associate professor; during that time, he also taught English
11 linguistics at the University of Hamburg. In 2005, he was at the
12 Psychology Department of the Max Planck Institute for Evolutionary
13 Anthropology in Leipzig before he moving to UCSB in November 2005.
14 Gries was a visiting professor at five LSA Linguistic Institutes, a
15 Visiting Chair (2013-2017) of the Centre for Corpus Approaches to
16 Social Science at Lancaster University, and the Leibniz Professor
17 (spring semester 2017) at the Leipzig University.
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60