

1
2
3
4 **Typological parameters of intralingual variability:**
5 **Grammatical analyticity versus syntheticity in varieties**
6 **of English**
7

8 BENEDIKT SZMRECSANYI
9 *Freiburg Institute for Advanced Studies*
10

11
12 ABSTRACT

13 Drawing on terminology, concepts, and ideas developed in quantitative morphological
14 typology, the present study takes an exclusive interest in the coding of grammatical
15 information. It offers a sweeping overview of intralingual variability in terms of
16 overt grammatical *analyticity* (the text frequency of free grammatical markers),
17 grammatical *syntheticity* (the text frequency of bound grammatical markers), and
18 *grammaticity* (the text frequency of grammatical markers, bound or free) in English.
19 The variational dimensions investigated include geography, text types, and real time.
20 Empirically, the study taps into a number of publicly accessible text corpora that
21 comprise a large number of different varieties of English. Results are interpreted in
22 terms of how speakers and writers seek to achieve communicative goals while
23 minimizing different types of complexity.

24 This study surveys language-internal variability and short-term diachronic change,
25 along dimensions that are familiar from the cross-linguistic typology of languages.
26 The terms *analytic* and *synthetic* have a long and venerable tradition in linguistics,
27 going back to the 19th century and August Wilhelm von Schlegel, who is usually
28 credited for coining the opposition (Schlegel, 1818). This is, alas, not the place to
29 even attempt to review the rich history of thought in this area (but see Schwegler,
30 1990:chapter 1 for an excellent overview). Suffice it to say that the terms “are used
31 in widely different meanings by different linguists” (Anttila, 1989:315), a
32 terminological confusion that requires, right at the outset, a concise definition
33 that guides the present study’s empirical argument (note that *grammaticity*, as a
34 derived notion, will be defined in a following section). This study is interested,
35 first, in the overt coding of *grammatical* information, which is why lexical
36 analyticity and syntheticity do not enter into consideration here. Second, our
37 definition is a strictly *formal* one (and not a semantic one) that broadly follows
38 Andrei Danchev’s notion that “formal analyticity ... implies that the various
39 meanings ... of a given language unit are carried by ... free morphemes, whereas
40

41
42 The author is grateful to Nicholas Smith for making available the CLAWS8-annotated versions of the
43 LOB and Brown corpora, and to Christian Mair for providing access to the Jamaican component of
44 ICE, even though the component is not yet officially released. The author also wishes to thank Peter
45 Auer, Douglas Biber, Stefan Th. Gries, Bernd Kortmann, and three anonymous referees for helpful
comments and suggestions. The usual disclaimers apply.

46 formal syntheticity is ... characterized by the presence of one bound morpheme”
 47 (Danchev, 1992:26). In this spirit—but avoiding reference to the *morpheme*
 48 construct, which is theoretically not unproblematic—we operationally define:

49
 50 *Formal grammatical analyticity*—comprises all those coding strategies in which
 51 grammatical information is conveyed by free grammatical markers, which we in
 52 turn define as *synsemantic* (cf. Marty, 1908) word tokens that have no independent
 53 lexical meaning.

54 *Formal grammatical syntheticity*—comprises all those coding strategies where
 55 grammatical information is signaled by *bound* grammatical markers.

56
 57 A few additional remarks are in order here. As for analyticity, this study equates
 58 synsemantic word tokens with *function* (also known as *structure* or *empty*) *words*,
 59 which are here defined as being members of closed word classes: conjunctions
 60 (e.g., *and*, *if*), determiners (e.g., *the*), pronouns (e.g., *he*), prepositions (e.g., *in*),
 61 infinitive markers (e.g., *to*), modal verbs (e.g., *can*, *will*), and negators (e.g., *not*).
 62 Note that this definition of analyticity and of what should count as a function word
 63 appears to be a fairly uncontroversial one and in accordance with standard
 64 reference works (for instance, Bussmann, Trauth, & Kazzazi, 1996:22, 471). As
 65 for our definition of syntheticity, we take bound grammatical markers to comprise
 66 verbal, nominal, and adjectival inflectional affixes (e.g., past tense *-ed*, plural *-s*,
 67 comparative *-er*, and so on), the genitive clitic (as in *Tom’s house*), as well as
 68 allomorphs including ablaut phenomena (e.g., past tense *sang*), *i*-mutation (e.g.,
 69 plural *men*), and other nonregular yet clearly bound grammatical markers. Our
 70 model of morphological analysis is thus, at base, an *item-and-process* model
 71 (Hockett, 1954:396) in which grammatically marked forms are thought of as
 72 deriving from simple forms via some sort of process—in our diction, via adding
 73 some sort of overt grammatical (but not necessarily segmentable) grammatical
 74 signal, be it a (regular) inflectional affix, a stem vowel change, or the like. What
 75 does not enter into our notion of syntheticity, however, is the “zero morpheme”
 76 construct (as in *they go-Ø*) postulated in some morphological approaches to deal
 77 with paradigmatic contrasts in finite verb forms. Note here that the present study is,
 78 in fact, going to be interested in null marking—but only to the extent that null
 79 marking serves as an *alternative* to non-null synthetic (and also analytic) marking,
 80 not as an *instantiation* of synthetic marking.

81 Having thus set the scene, what is the main objective in this article? In cross-
 82 linguistic morphological typology, languages are classified as rather analytic (for
 83 instance, modern Romance languages) or as rather synthetic (for example,
 84 Classical Latin). English is frequently cited as the textbook example of a
 85 language that has developed from a synthetic language into an analytic one.
 86 Consider the following, classical Schlegel quote:

87
 88 En Europe les langues dérivées du latin, et *l’anglais*, ont une grammaire tout
 89 analytique ... synthétiques dans leur origine ... elles penchent fortement vers les
 90 formes analytiques.¹ (Schlegel, 1846:161; emphasis mine)

91 Against this backdrop, the present study seeks to demonstrate that Modern English is
 92 not as monolithic, or monolithically analytic, as the preceding quote and indeed much
 93 of the cross-linguistic morphological literature would seem to suggest. Instead, we
 94 will see that variability in analyticity and syntheticity is endemic, surprisingly so,
 95 even among closely related dialects and varieties of the same language, Modern
 96 English. The task before us is to marry Schlegel’s old idea about the syntheticity-
 97 analyticity continuum to state-of-the-art corpus-linguistic techniques along the
 98 lines of Gries (2006) to explore three language-external parameters of intralingual
 99 variability: geography, text type, and short-term diachrony.

100 A comment on this study’s general methodological orientation seems
 101 appropriate at this point. Though interested in variation and variability, the
 102 present study does not adopt a strictly variationist approach in the sense of, for
 103 example, Labov (1966, 1972). Whereas it is certainly possible, in many cases, to
 104 define a linguistic variable that has an analytic variant and a synthetic variant
 105 (e.g., the analytic *of*-genitive versus the synthetic *s*-genitive, analytic adjective
 106 comparison vs. synthetic adjective comparison, and so on); in the majority of
 107 instances, a particular analytic or synthetic pattern will not have a neatly
 108 definable alternative variant. For example, synthetic marking of plurality (e.g.,
 109 *many horses*) does not have an analytic alternative that would convey exactly or
 110 even roughly the same meaning. As the subsequent frequency analyses will
 111 show, however, the alternative to analytic and synthetic marking is, in many
 112 contexts, no grammatical marking at all; the empirical question being whether
 113 analytic or synthetic marking is more likely to be substituted by zero.

114
 115
 116 INDICES

117
 118 In a seminal (1960) paper entitled “A Quantitative Approach to the Morphological
 119 Typology of Language,” Joseph Greenberg demonstrated that *prima facie* abstract
 120 typological notions are amenable to sufficiently precise numerical measurements
 121 by calculating a number of indices on the basis of naturalistic texts. Greenberg
 122 defined (i) an index of synthesis, (ii) of agglutination, (iii) a compounding index,
 123 (iv) a derivational index, (v) a gross inflectional index, (vi) a prefixial index,
 124 (vii) a suffixial index, (viii) an isolational index, (ix) a pure inflectional
 125 index, and (x) a concordial index (Greenberg, 1960:187). So, for instance,
 126 Greenberg defined the gross inflectional index as the number of inflectional
 127 morphemes (nonconcordial or concordial) in the analyst’s sample divided by the
 128 total number of words in the sample (Greenberg, 1960:186–187).²

129 This study utilizes Greenberg’s method in revised form. First, as for inflection it
 130 calculates *syntheticity indices* in a slightly different fashion. What is measured, in a
 131 given textual sample, is not the number of inflectional morphemes per sample
 132 (which is what Greenberg’s original gross inflectional index measures), but the
 133 number of words in a sample that bear at least one bound grammatical marker. Note
 134 that these are not necessarily two different ways of saying the same thing, as—
 135 depending on one’s analytical framework—the form *walks* (as in *he walks the dog*)

136 can be analyzed as containing two grammatical morphemes, {nonpast} and {third-
 137 person singular}. In our approach, the form *walks* contains exactly one grammatical
 138 marker, *-s*, which may have more than one meaning. Notice here that except for
 139 some rare genitive plural forms (such as *the oxen's legs*), English has virtually no
 140 word forms that exhibit more than one segmentable bound grammatical marker.

141 Second, what is notably absent from Greenberg's index portfolio is an
 142 analyticity index. In an attempt to remedy this omission, Kasevič & Jachontov
 143 (1982:37) (cited in Kempgen & Lehfeldt, 2004:1237) suggested an "index of
 144 analyticity," which relates the number of synsemantic words in a given text to
 145 the total number of words in that text (cf. Kelemen, 1970:62 for a similar
 146 proposal). This is how the present study calculates its *analyticity index*.

147 In addition to the syntheticity and analyticity indices, we calculate *grammaticity*
 148 *indices* that measure the number of grammatical markers, free *or* bound, in a given
 149 sample. Numerically, the grammaticity index equals the sum of the former two
 150 indices. Grammaticity is equivalent to grammar minus word order in that the
 151 notion comprises all explicit grammatical markers, but not word order. To
 152 summarize, the present study is concerned with calculating three different indices:

153

- 154 1. The *analyticity index* (henceforth AI): the ratio of the number of free grammatical
 155 markers in a sample (F) to the total number of words in the sample (W),
 156 normalized to a sample size of 1,000 tokens. Hence: $AI = F/W \times 1,000$.
- 157 2. The *syntheticity index* (henceforth SI): the ratio of the number of words in a
 158 sample that bear a bound grammatical marker (B) to the total number of words
 159 in the sample (W), normalized to a sample size of 1,000 tokens. Hence: $SI =$
 160 $B/W \times 1,000$.
- 161 3. The *grammaticity index* (henceforth GI): the ratio of the total number of grammatical
 162 markers ($B + F$) in a text to the total number of words (W) in the sample,
 163 normalized to a sample size of 1,000 tokens. Hence: $GI = (B + F)/W \times 1,000$.

164 All three indices have a lower bound of zero. The syntheticity and the analyticity
 165 index have an upper bound of 1,000 index points, whereas the grammaticity index
 166 has an upper bound of 2,000 index points.

167

168

169 THE LINK TO LANGUAGE COMPLEXITY

170

171 A survey of the literature reveals that there is a customary nexus between
 172 analyticity, syntheticity, and *language complexity* (note, though, that this nexus
 173 is not always backed up by hard empirical evidence, especially when it comes to
 174 processing complexity). Be that as it may, the present study is in line with a
 175 number of orthodox interpretational patterns, which can be summed up as follows.

176 Wilhelm von Humboldt was one of the first to claim that analyticity increases
 177 explicitness and transparency while easing comprehension difficulty (Humboldt,
 178 1836:284–285). Syntheticity, on the other hand, is often viewed as increasing
 179 speaker/writer output economy and expressivity (cf. Danchev, 1992:36) by virtue
 180 of the fact that synthetic marking (in English, typically affixation) is the more

181 compact and economical coding option vis-à-vis analytic marking. Consider the
 182 alternation between the synthetic *s*-genitive (as in *the president's speech*) and
 183 the analytic *of*-genitive (as in *the speech of the president*). The synthetic option is
 184 more output-economical, because the genitive marker is a clitic and not a full-
 185 blown preposition, and because the possessed NP lacks a determiner. But, by the
 186 same token, the analytic option is the more explicit and arguably the more
 187 transparent one, by virtue of the fact that more material is used for grammatical coding.

188 Thus, syntheticity is more output-economical than analyticity is because synthetic
 189 markers are typically more compact than analytic markers. As for grammaticity, this
 190 study equates—following the argument in Szmrecsanyi & Kortmann (2009)—
 191 increased text frequencies of grammatical markers (synthetic *or* analytic) with
 192 “repetition of information” (Trudgill, 2009: 314), which increases marking
 193 redundancy and hence decreases overall speaker/writer output economy, because
 194 more overt grammatical information is being explicitly coded, be it synthetically or
 195 analytically. This is why we consider less grammaticity to be more output-
 196 economical than more grammaticity, and here it does not matter if more
 197 grammaticity comes about through more analyticity or syntheticity, because zero
 198 marking is *always* more output-economical than explicit marking is. On the other
 199 hand, however, increased grammaticity, that is, increased overt redundancy, can
 200 be seen (for instance, Bisang, 2009) as easing hearer/reader pragmatic inference
 201 complexity, because less is left for the reader/hearer to pragmatically infer from the
 202 context. So, the basic idea, advocated in, for example, Bisang (2009), is that
 203 there is a trade-off between competing motivations such that higher levels of
 204 grammaticity are comprehension-economical with regard to the hearer/reader
 205 whereas lower levels of grammaticity are output-economical with regard to the
 206 speaker/writer. Our interpretational approach can thus be summarized as follows:

- 207
- 208 1. Increased analyticity increases explicitness and transparency and decreases
 209 hearer/reader comprehension complexity.
 - 210 2. Increased syntheticity increases speaker/writer output economy vis-à-vis analytic
 211 marking, by virtue of being the more compact coding option.
 - 212 3. Increased grammaticity (i) increases redundancy, thus (ii) decreasing overall
 213 speaker/writer output economy, because more grammatical information is
 214 subject to overt coding. Redundancies such as these, however, (iii) reduce
 215 hearer/reader pragmatic inference complexity.
- 216
217

218 METHOD AND DATA

219
220 *Data*

221
222 The present study draws on a fairly wide array of publicly accessible text corpora:

223
224 *The British National Corpus* (BNC World Edition). This data base contains
 225 approximately 90 million words of written standard British English (henceforth:

226 BrE) and 10 million words of spoken standard BrE. Containing over 4,000 individual
 227 texts, the corpus samples 70 different registers (24 spoken, 46 written; for instance,
 228 S_speech_scripted vs. S_speech_unscripted, W_fict_drama vs. W_fict_poetry)
 229 at the highest level of granularity, which boil down to 34 macro registers
 230 (16 spoken, 18 written; for instance, S_speech and W_fict). The corpus is fully
 231 part-of-speech (henceforth: POS) annotated using the CLAWS5 tag set (Aston &
 232 Burnard, 1998).

232 *The Brown family of corpora* (Brown, LOB, Frown, F-LOB). These are four matching
 233 text corpora sampling 1960s American English (henceforth: AmE) (Brown) (see
 234 Francis & Kučera, 1982), 1960s BrE (LOB) (see Johansson & Hofland, 1989),
 235 1990s AmE (Frown) (see Hinrichs, Waibel, & Smith, 2007), and 1990s BrE
 236 (F-LOB) (see Hinrichs et al., 2007). The four corpora have (roughly) the same
 237 design, each spanning one million words (500 texts of approximately 2,000 words
 238 each) and sampling 15 written micro registers falling into four macro registers
 239 (Press, General Prose, Learned Writing, Fiction) and two major text categories
 240 (informative vs. imaginative prose). The corpora are fully POS-annotated with the
 241 CLAWS8 tag set.

242 *Switchboard*. This corpus samples AmE telephone conversations. The version that is
 243 used here stems from the second release of the *American National Corpus*.³
 244 It contains approximately three million words and is POS-annotated with the
 245 Hepple tag set. *Switchboard* serves to represent standard spoken AmE.

246 *The Freiburg Corpus of English Dialects* (FRED). FRED contains transcribed oral
 247 history interviews, the bulk of which were recorded in the 1970s and 1980s.
 248 Speakers are typically non-mobile old rural males. FRED yields three levels of
 249 areal granularity: 9 dialect areas, 38 counties (in pre-1974 boundaries), and 163
 250 locations (see Hernández, 2006; Szmrecsanyi & Hernández, 2007). This study
 251 explores variation between six dialects sampled in FRED: the dialects spoken in
 252 the counties of Somerset (southwestern England), Kent (southeastern England),
 253 Shropshire (English Midlands), Lancashire (northern England), Glamorgan
 254 (Wales), and Sutherland (Scottish Lowlands).⁴

255 *The International Corpus of English (ICE)*. The following ICE subcorpora are
 256 analyzed: ICE-IRE (Irish E), ICE-PHI (Philippine E), ICE-HK (Hong Kong E),
 257 ICE-SG (Singapore E), ICE-IN (Indian E), ICE-NZ (New Zealand E), ICE-JAM
 258 (Jamaican E), and ICE-EA (East African E, i.e., Kenyan and Tanzanian E).⁵ The
 259 subcorpora typically contain 500 texts (300 spoken, 200 written), with every
 260 individual text spanning approximately 2,000 words (Greenbaum, 1996). Of the
 261 many registers sampled in ICE, this study explores the spoken-conversational
 262 material (section 11a).

263 *Method*

265 To obtain quantitative results, the present study exploits POS annotation, where
 266 tokens in a corpus are tagged for their word class (this includes information on
 267 whether nouns, verbs, adjectives, and certain pronouns carry inflections). In the
 268 case of those corpora that are POS-annotated in the first place (i.e., the BNC,
 269 *Switchboard*, and the *Brown* family of corpora), the findings derive from an
 270 exhaustive analysis of *all* the material sampled in the corpora.

271 In the case of those corpora in the present study’s portfolio that are not POS-
 272 annotated a priori (essentially, FRED and the ICE subcorpora), an algorithm
 273 selected 1,000 random decontextualized tokens (i.e., words) per variety studied.
 274 Subsequently, these tokens were annotated manually for their part of speech
 275 using the BNC (CLAWS5) tag set with a minor extension (as in the CLAWS8
 276 tag set, the primary verbs *be*, *do*, and *have* were explicitly annotated for whether
 277 they occurred in auxiliary function by prefixing the character ‘A’ to the
 278 CLAWS5 tag; note that in the analysis of the BNC itself, primary verbs were
 279 automatically disambiguated contextually for auxiliary or main verb usage).⁶

280 Given the definition of analyticity and syntheticity detailed in the first section,
 281 POS tags (or rather the tokens annotated with POS tags) were subsequently placed
 282 into four categories: (i) purely lexical tags, such as singular nouns (which are
 283 uninteresting to the present study), (ii) synthetic tags (essentially all tokens that,
 284 following Vennemann, 1982:330, show affixation or mutation to indicate
 285 grammatical information), (iii) analytic tags (at base, function words), and (iv) a
 286 small number of simultaneously synthetic and analytic tags (inflected auxiliary
 287 verbs and reflexive pronouns in their plural form). The exact tag/token-to-
 288 category matches can be seen in Tables 1⁷ and 2, which categorize analytic tags
 289 into 11 broad component categories and synthetic tags into 4 broad component
 290 categories.⁸

291 Finally, a retrieval script written in the programming language Perl
 292 automatically established the text frequencies of the relevant POS tags (or POS-
 293 tag categories) in the data set. These text frequencies served as the empirical
 294 basis for calculating the indices.⁹

295
 296 GEOGRAPHIC VARIABILITY IN WORLD ENGLISHES
 297

298 This section explores analyticity, syntheticity, and grammaticity variability in 16
 299 geographic varieties of English, comprising 10 geographic L1 varieties of
 300 English and 6 non-native, indigenized L2 (or ESL) varieties of English, all of
 301 which are spoken. The L1 varieties include the traditional dialects spoken in
 302 Glamorgan, Kent, Lancashire, Shropshire, Somerset, and Sutherland; Irish E;
 303 New Zealand E; standard (conversational) BrE; and standard spoken AmE.
 304 The L2 varieties are East African E, Indian E, Hong Kong E, Philippine E,
 305 Singapore E, and Jamaican E.¹⁰ Crucially, in controlling for text type as far as
 306 possible (notice that the data subject to analysis in this section are all
 307 spontaneous-spoken), the variability subject to analysis in this section can be
 308 considered genuinely geographic.

309
 310 *Analyticity and syntheticity in World Englishes*

311
 312 For every one of the 16 geographic varieties under investigation, the scatterplot in
 313 Figure 1 plots analyticity indices against syntheticity indices in a two-dimensional
 314 plane, differentiating visually between native L1 varieties and non-native,
 315 indigenized L2 varieties. It is evident that variability in World Englishes is

TABLE 1. *Eleven broad component categories (as defined through POS tags and/or word tokens) loading on the analyticity index in (i) the modified BNC tag set (CLAWS5e) used for manual annotation, (ii) the original BNC tag set (CLAWS5), (iii) the Brown family tag set (CLAWS 8), and (iv) the Switchboard (Hepple) tag set*

		CLAWS5e tags (manual coding)	CLAWS5 tags (BNC)	CLAWS8 tags (Brown family)	Hepple tags (Switchboard)
1	conjunctions, subjunctions, prepositions	CJ*, PRF, PRP	CJ*, PRF, PRP	C*, WPR, I*	CC, IN
2	determiners, articles, <i>wh</i> - words	D*, AT0, AVQ, PNQ	D*, AT0, AVQ, PNQ	APPGE, AT*, D*, RGQ, RGQV, RRQ, RRQV	DT, W*
3	existential <i>there</i>	EX0	EX0	EX	EX
4	pronouns	PNI, PNP, PNX	PNI, PNP, PNX	P*	PP, PRP, PRPS, PRPRS
5	<i>more, most</i>	<i>more, most</i>	<i>more, most</i>	<i>more, most</i>	<i>more, most</i>
6	infinitive marker <i>to</i>	TO0	TO0	TO	TO + VB
7	modals	VM0	VM0	VM*	MD
8	negator <i>not, n't</i>	<i>not, n't</i>	<i>not, n't</i>	<i>not, n't</i>	<i>not, n't</i>
9	auxiliary BE ^a	AVB*	VB* + V*, VB* + * + V*, VB* + XX0	VAB* + V*, VAB* + * + V*, VAB* + XX0	TO BE + V, TO BE + * + V, TO BE + <i>not</i> , TO BE + *n't
10	auxiliary DO ^a	AVD*	VD* + V*, VD* + * + V*, VD* + XX0	VAD* + V*, VAD* + * + V*, VAD* + XX0	TO DO + V, TO DO + * + V, TO DO + <i>not</i> , TO DO + *n't
11	auxiliary HAVE ^a	AVH*	VH* + V*, VH* + * + V*, VH* + XX0	VAH* + V*, VAH* + * + V*, VAH* + XX0	TO HAVE + V, TO HAVE + * + V, TO HAVE + <i>not</i> , TO HAVE + *n't

^aCategory members may also load on the syntheticity index.

considerable. In the analyticity dimension (vertical axis), values range from 403 analytic markers per 1,000 words of running text (Hong Kong E) all the way to 531 markers per 1,000 words (Sutherland English). In terms of syntheticity (horizontal axis), values span the range between 111 synthetic markers per 1,000 words (Singapore E) and 185 markers per 1,000 words (Shropshire E). The two standard varieties, standard AmE and BrE, cover the middle ground. We shall now discuss some general tendencies in the data set, not all of which are significant, due to the comparatively low number of observations (such tendencies will be reported anyway, with the proviso that they are tentative and await corroboration in future research).

According to Peter Trudgill (2009), the distinction between *high-contact* and *low-contact varieties* of English is what amounts to “the true typological split”

TABLE 2. Five broad component categories (as defined through POS tags and/or word tokens) loading on the syntheticity index in (i) the modified BNC tag set (CLAWS5e) used for manual annotation, (ii) the original BNC tag set (CLAWS5), (iii) the Brown family tag set (CLAWS 8), and (iv) the Switchboard (Hepple) tag set

		CLAWS5e tags (manual coding)	CLAWS5 tags (BNC)	CLAWS8 tags (Brown family)	Hepple tags (Switchboard)
12	<i>s</i> -genitive	POS	POS	GE, MCGE	POS
13	comparative and superlative adjectives	AJC, AJS	AJC, AJS	JJR, JJT	JJR, JJS, JJSS ^a
14a	plural nouns	NN2	NN2	NN2,>NNL2, NNO2, NNT2, NNU2, NP2, NPD2, NPM2	NNPS, NNS, NPS
14b	plural reflexive pronouns ^b	PNX+ word token ending in <i>*ves</i>	PNX+ word token ending in <i>*ves</i>	PPX2	<i>*selves</i>
15	inflected verbs ^b	VVD, VVG, VVN, VVZ, (A) VBB, ^c (A) VBD, (A) VBG, (A) VBN, (A) VBZ, (A) VDD, (A) VDG, (A) VDN, (A) VDZ, (A) VHD, (A) VHN, (A) VHZ	VVD, VVG, VVN, VVZ, VBB, ^c VBD, VBG, VBN, VBZ, VDD, VDG, VDN, VDZ, VHD, VHG, VHN, VHZ	VVBDR, VVBDZ, VVBG, VVBM, VVBN, VVBR, VVBZ, VVDD, VVDG, VVDN, VVDZ, VVHD, VVHG, VVHN, VVHZ, VABDR, VABDZ, VABG, VABM, VABN, VABR, VABZ, VADD, VADG, VADN, VDZ, VAHD, VAHG, VAHN, VAHZ, VVD, VVG, VVGK, VVN, VVNK, VVZ	VBD, VBG, VBN, VBZ

^aConsidered if and only if not preceded by *more* or *most*.

^bCategory members may also load on the analyticity index.

^cConsidered if and only if the form of the verb is not *be*.

(Trudgill, 2009:315) among varieties of English. Trudgill’s argument boils down to the purportedly “lousy language-learning abilities of the human adult” (Trudgill, 2001:372). The idea, in a nutshell, is that contact implicates adult language learning, which in turn implicates simplification. The resulting simplicity in different domains (see Trudgill, 2009, for an overview) is what should set apart high-contact from low-contact varieties as synchronic groups. In terms of the present data set, high-contact varieties thus comprise:

Non-native, indigenized L2 varieties: East African E, Indian E, Hong Kong E, Philippine E, Singapore E, and Jamaican E;

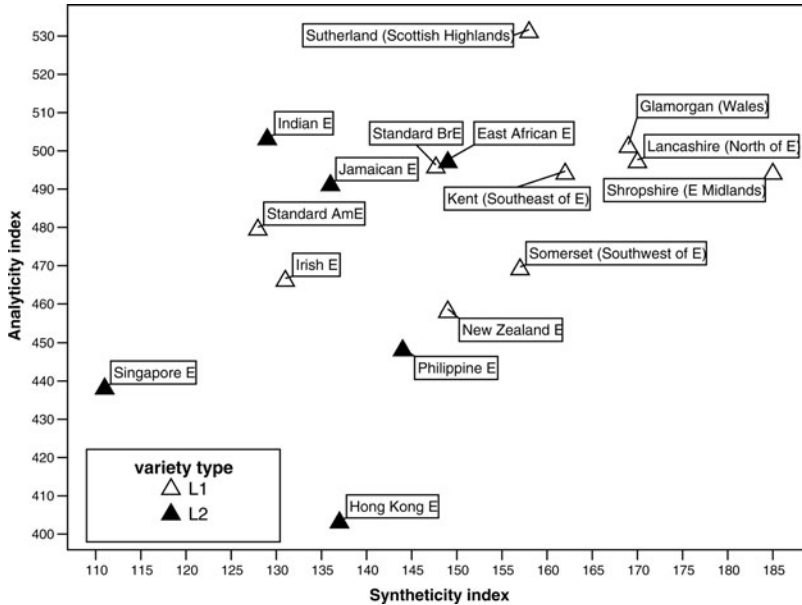


FIGURE 1. Geographic varieties: analyticity by syntheticity (in index points, ptw).

Transplanted L1 Englishes (or *colonial varieties*; cf. Mesthrie, 2006a:382): New Zealand E and standard spoken AmE;

Language-shift Englishes: varieties “that develop when English replaces the erstwhile primary language(s) of a community” and that have “adult and child L1 and L2 speakers forming one speech community” (Mesthrie, 2006a:383). The present study also includes what might be called *shifted varieties*, which are varieties that used to be genuine language-shift varieties in the past 500 years or so but which no longer have significant numbers of L2 speakers any more: Irish E, Welsh E (Glamorgan), Scottish Highlands E (Sutherland);

Standard varieties, such as standard BrE and standard AmE, the genesis of which, according to Trudgill (2009), *always* implicates a high degree of dialect contact.

Varieties that do not fall into one of the preceding categories are considered low-contact L1 dialects of English, that is, traditional nontransplanted regional dialects that are “long-established mother tongue varieties” (Trudgill, 2009:320); thus, in terms of the data set analyzed here, the traditional dialects spoken in Kent, Lancashire, Shropshire, and Somerset.

Observe, first, that low-contact varieties are significantly ($p = .006$)¹¹ more synthetic than high-contact varieties (mean SI high-contact: 141; mean SI low-contact: 169) (cf. Szmrecsanyi & Kortmann, 2009).¹² Second, L1 varieties exhibit significantly ($p = .023$) more syntheticity (mean SI: 156) than L2 varieties (mean SI: 134). Among L2 varieties, there is a striking gap between Southeast Asian L2 varieties (Singapore E, Philippine E, Hong Kong E) and non-Southeast Asian L2 varieties (Indian E, Jamaican E, East African E)

451 inasmuch as the former exhibit significantly ($p = .009$) less analyticity than do the
 452 latter. There is also a tendency for Southeast Asian varieties to exhibit less
 453 syntheticity than other L2 varieties. We will return to this issue later.

454 Third, the standard view in the literature (on English or other European
 455 languages) is that there is a historical trade-off between syntheticity and
 456 analyticity. English, for instance, is said to have compensated for the loss of
 457 synthetic marking by adding analytic marking (cf. the discussion and quotes in
 458 the first section). What is interesting about Figure 1 is that on the synchronic-
 459 geographic plane, there is no such thing as a trade-off between analyticity and
 460 syntheticity; on the contrary, there seems to be a positive ($r = .47$), though
 461 statistically marginally insignificant ($p = .066$), correlation between analyticity
 462 and syntheticity. In short, varieties of English that exhibit greater syntheticity
 463 tend to also exhibit a high degree of analyticity, whereas varieties that display
 464 little syntheticity will typically also have little analyticity. The crucial variable
 465 thus seems to be grammaticity, to which we turn next.

466
 467 *Grammaticity*

468 Table 3 provides GI scores for every one of the 16 geographic varieties under study
 469 in this section, along with z scores as indices of variability between varieties. The
 470 variety exhibiting the lowest level of grammaticity is Hong Kong E: 539
 471 grammatical markers per 1,000 words of running text. Its z score of -2.0
 472 indicates that Hong Kong E's GI is 2.0 cross-variety standard deviations less
 473 than the mean of all varieties under study (in the data set, the standard deviation
 474 for grammaticity is 43.4 index points; the mean index value for grammaticity is
 475 626.9). By contrast, Sutherland E has a GI of 689.0, which is 1.4 standard
 476
 477

478 TABLE 3. *Grammaticity in geographic varieties of English*

	Grammaticity index	
	Mean	z score
482 Hong Kong E	539	-2.0
483 Singapore E	549	-1.8
484 Philippine E	592	-.8
485 Irish E	598	-.7
486 New Zealand E	607	-.5
Standard AmE	607	-.4
487 Somerset (southwestern England)	626	.0
488 Jamaican E	627	.0
489 Indian E	632	.1
Standard BrE	643	.4
490 East African E	647	.5
491 Kent (southeastern England)	657	.7
492 Lancashire (northern England)	667	.9
Glamorgan (Wales)	669	1.0
493 Shropshire (English Midlands)	680	1.2
494 Sutherland (Scottish Highlands)	689	1.4

495

496 deviations greater than the mean of all varieties. In summary, starting at the top of
 497 Table 3, Hong Kong E through Somerset E have lower-than-average grammaticity;
 498 Jamaican E through Sutherland E exhibit higher-than-average grammaticity.

499 Again, some tendencies in the data set should be pointed out. For one thing,
 500 British varieties of English tend to exhibit more grammaticity (mean GI: 654)
 501 than other varieties, whereas L2 varieties have less-than-average grammaticity
 502 (mean GI: 598), with transplanted L1 varieties taking the middle road (mean GI:
 503 607); these differences are significant at $p = .032$, according to a one-way
 504 analysis of variance. In a similar vein, L1 varieties tend to exhibit more
 505 grammaticity (mean GI: 644) than L2 varieties do (mean GI: 598; $p = .032$).
 506 Again, we find the three Southeast Asian L2 varieties (Singapore E, Philippine
 507 E, Hong Kong E) at the bottom of Table 3. With a mean GI of merely 560,
 508 these differ significantly ($p = .001$) from the other varieties in the sample in that
 509 they seem to avoid grammatical marking. At this point, it is instructive to
 510 examine some authentic text samples: (1) is a conversational snippet taken from
 511 the Hong Kong E data, and (2) exemplifies conversational Singapore E. Sites
 512 where grammatical markers could appear are marked by \emptyset .

513

514

- 515 1. ...the Putonghua we we speak, uhm, include- \emptyset just part of the Beijing dialect.
 516 But in Beijing most people especially the, uhm people who are less educated
 517 they speak \emptyset Beijing dialect which is really difficult for us to understand. Even
 518 I myself I finished \emptyset advance course I, can't hear what they say. (ICE-HK
 519 S1A-002)
- 520 2. ...and he actually went you know to Robinson and bought him two shirt- \emptyset and
 521 one tie ... Ah because he lost the bet mah because he say- \emptyset that if Wei Ho
 522 change- \emptyset he will do that for him ... But when he came in on Monday that day
 523 we almost die- \emptyset laughing uh and Kang Heng also lost the bet to us ... Kang
 524 Heng say- \emptyset he won't change If he change- \emptyset he give- \emptyset us ten dollars. (ICE-
 525 SG S1A-013)

525

526

527 Consider, now, Hong Kong E. In (1), we find two sites (*they speak \emptyset Beijing*
 528 *dialect, I finished \emptyset advance course*) where an article could be employed, and
 529 one site (*the Putonghua ... include- \emptyset just part of the Beijing dialect*) where
 530 many speakers of L1 standard varieties would employ a verbal inflection.
 531 Likewise, in the relatively short snippet in (2), we find seven unmarked nouns or
 532 verbs (*two shirt- \emptyset , he say- \emptyset , Wei Ho change- \emptyset , we almost die- \emptyset , Kang Heng*
 533 *say- \emptyset , If he change- \emptyset he give- \emptyset us*) where there could be inflectional forms.
 534 The generalization seems to be that in Southeast Asian L2 varieties of English in
 535 particular, speakers do not substitute, say, synthetic grammatical markers by
 536 purportedly more transparent and explicit analytic markers. Instead, they opt for
 537 less overall grammaticity, avoiding overt marking—in the spirit of the motto “if
 538 it can be deleted, it will be deleted” (Mesthrie, 2006b:142). This strategy is even
 539 more output-economical than synthetic marking is, yet it arguably incurs
 540 pragmatic complexities on the part of the hearer (cf. Bisang, 2009).

The sources of geographic analyticity/syntheticity variability

Let us now identify those individual grammatical markers and/or marker categories that are most strongly involved in geographic analyticity/syntheticity variability. This means that we will deconstruct the indices considered so far, elucidating which of the 15 component categories detailed in Tables 1 and 2 are subject to significant variability. We begin by exploring analytic markers. Correlating, in our sample of $N=16$ varieties, each of the component categories with the analyticity index yields a statistically significant ($p < .003$)¹³ Pearson correlation coefficient for articles, determiners other than articles, and *wh*-words. These correlate strongly ($r=.83$) with increased AI levels. Subsequent independent samples *t* tests show that the text frequency of such markers is also what sets Southeast Asian L2 varieties apart from other varieties of English ($p = .007$). Within this category it is especially articles that show robust variability (see example (1)). As for syntheticity, we obtain an even stronger correlation ($r=.90$) between increased SI levels and inflected verbs (especially third-person singular and past tense forms). Once again, independent samples *t* tests indicate that inflected verbs are likewise highly involved in the overall divide between Southeast Asian L2 varieties and other varieties ($p=.029$). Recall that Singapore E is the least synthetic variety in our sample according to Figure 1, and as we have seen in example (2), this variety exhibits a strong preference for unmarked verb forms.

Interim summary

In all three relevant dimensions—analyticity, syntheticity, and grammaticity—geographic varieties of English are subject to substantial variability. This section has offered the following generalization. Low-contact varieties are more synthetic than high-contact varieties.¹⁴ Thus, in the data set subject to analysis here, low-contact communities emphasize output economy whereas high-contact speaker communities put a premium on explicitness and transparency. Furthermore, L1 varieties exhibit more grammaticity than L2 varieties do, and Southeast Asian L2 varieties in particular are substantially less explicit grammatically than are other varieties. Hence, Southeast Asian L2 varieties are more economical than other varieties when it comes to grammaticity (note that this seems to be true for Southeast Asian languages in general, according to Bisang, 2009). Next, exploring which grammatical markers are specifically involved in this kind of variability, we have seen that determiners, articles, and *wh*-words—and within this category, articles more than anything else—are loading high on the analyticity index. As for syntheticity, it is mainly text frequencies of verbal inflections (or the lack thereof) that are most strongly implicated in the observable variability. Last but not least, analyticity does not seem to trade off against syntheticity such that reduced syntheticity would imply increased analyticity or vice versa. Instead, there is not only a binary choice between analyticity and syntheticity, but also a third option—zero—which is often the preferred one.

586 TEXT TYPE VARIABILITY

587
 588 This portion of the article investigates text type (or genre) variability in standard
 589 BrE, drawing on the BNC as the primary corpus database. Recall that the point
 590 estimate for spoken Standard BrE in the previous section (cf. Figure 1) was
 591 based on the conversational part of the BNC. What picture would emerge if we
 592 explored variability between this genre and the many other registers sampled in
 593 the BNC? It is to this task that we next turn.

594

595 *Analyticity and syntheticity*

596

597 We begin by looking at analyticity-syntheticity variability. The BNC in its entirety
 598 yields an AI of 440.2 and an SI of 176.9, but, needless to say, there is a good deal
 599 of variability in the corpus. Table 4 reports some summary measures of this
 600 variability, at three granularity levels: individual texts, micro registers (for
 601 instance, *unscripted speech* vs. *scripted speech*, *drama fiction* vs. *prose fiction*),
 602 and macro registers (for instance, *speech* vs. *fiction*). For a first impression of
 603 this variability, observe that the least analytic individual text (text G2A, a
 604 collection of estate agents' property details) in the BNC only exhibits 228.5
 605 analytic markers per 1,000 words of running text, whereas the most analytic text,
 606 J98 (a Herts County Council committee meeting), is on record with an AI of
 607 570.8. As for the statistical dispersion around the mean, notice that the standard
 608 deviations are in the double-digit range and thus quite sizable, even at the level
 609 of macro registers.

610 In what follows, let us have a closer look at variability within and between macro
 611 registers. Table 5 lists the BNC's macro registers, along with point estimates for AI/
 612 SI, the standard deviation associated with each such point estimate (a measure of
 613 dispersion *within* macro registers), and z scores. So, for instance, broadcasts
 614 (S_brdbcast) have a mean AI of 472.4; the register-internal standard deviation
 615 associated with that mean is 33.5 (meaning that the 75 broadcast texts in the
 616 BNC deviate, on average, by 33.5 points from the preceding mean); and the z
 617 score is .4, which means that S_brdbcast's AI is .4 cross-macro register standard
 618 deviations greater than the mean of all macro registers' AI. A survey of the
 619 standard deviations in Table 5 reveals that there are more or less monolithic

620

621 TABLE 4. *Some summary measures (N of objects, minimum value, maximum value, standard*
 622 *deviation) for corpus-internal variability in the BNC, three levels of granularity: individual*
 623 *texts, micro registers, macro registers. Overall mean values: AI 440.2, SI 176.9*

	N	Analyticity index			Syntheticity index		
		Min	Max	SD	Min	Max	SD
Text types	4,052	228.5	570.8	47.0	46.9	271.9	21.9
Micro registers	70	358.0	548.4	39.9	127.3	197.1	18.4
Macro registers	34	378.8	548.4	39.2	131.9	192.5	17.9

630

TABLE 5. *Summary measures for variability between and within macro registers in the BNC: N of individual texts, point estimates for AI and SI (in index points, ptw), standard deviation, and z score*

	N	Analyticity index			Syntheticity index		
		Mean	SD	z score	Mean	SD	z score
S_brdcast	75	472.4	33.5	.4	160.4	18.0	-.1
S_classroom	58	485.0	22.5	.7	140.2	15.2	-1.2
S_consult	128	484.3	32.7	.7	134.2	20.3	-1.5
S_conv	153	495.7	24.8	1.0	147.7	13.7	-.8
S_courtroom	13	499.0	18.8	1.0	150.6	11.1	-.6
S_demonstratn	6	512.0	21.7	1.4	132.1	14.8	-1.6
S_interview	132	504.2	22.2	1.2	153.1	16.1	-.5
S_lect	31	482.1	27.1	.6	150.6	14.4	-.6
S_meeting	132	495.0	23.0	.9	146.3	11.6	-.8
S_parliament	6	476.5	19.1	.5	150.4	8.5	-.6
S_pub-debate	16	478.7	21.5	.5	132.0	8.6	-1.6
S_sermon	16	548.4	11.3	2.3	167.5	12.7	.3
S_speech	77	481.7	32.1	.6	149.0	16.5	-.7
S_sportslive	4	446.9	25.8	-.3	153.5	8.9	-.4
S_tutorial	18	491.7	31.8	.8	155.2	14.5	-.3
S_unclassified	44	486.3	57.3	.7	144.2	18.3	-.9
W_ac	500	437.1	32.3	-.5	179.5	17.3	1.0
W_admin	12	437.6	31.0	-.5	169.8	22.2	.5
W_advert	60	378.8	43.4	-2.0	158.5	22.4	-.2
W_biography	100	453.8	35.3	-.1	178.7	14.6	1.0
W_commerce	112	430.2	24.9	-.7	182.5	14.3	1.2
W_email	7	411.3	7.8	-1.2	148.4	4.0	-.7
W_essay	11	455.3	27.3	-.1	189.4	17.3	1.6
W_fict	464	480.5	25.2	.6	187.5	13.3	1.5
W_hansard	4	466.3	11.3	.2	154.4	5.8	-.3
W_institut-doc	43	407.7	26.4	-1.3	192.5	20.8	1.7
W_instructional	15	416.8	24.8	-1.1	162.3	15.4	.1
W_letters	17	444.3	25.9	-.4	145.0	38.5	-.9
W_misc	500	421.2	39.2	-.9	177.2	19.3	.9
W_news	32	395.8	13.2	-1.6	191.7	13.0	1.7
W_newsp	486	399.4	25.0	-1.5	182.1	17.0	1.2
W_non-ac	534	419.3	46.5	-1.0	181.0	14.9	1.1
W_pop	211	414.3	20.8	-1.1	172.3	12.9	.6
W_religion	35	475.7	22.4	.4	168.1	15.2	.4

genres. The seven *e-mail* texts sampled in the BNC are comparatively homogeneous with regard to the indices at hand, whereas the *unclassified* material is extraordinarily heterogeneous.

Register-internal variability aside, Figure 2 depicts the variability *between* macro registers by plotting AI/SI point estimates on a two-dimensional plane. Among the mass of data points displayed, a closer look at the extreme cases along the two dimensions in the diagram is instructive. In the syntheticity dimension, with SIs beyond 190, *institutional documents* and *news* are the most synthetic genres in the BNC, which is another way of saying that for these text types, the pressure for output economy is the strongest. At the other end of the

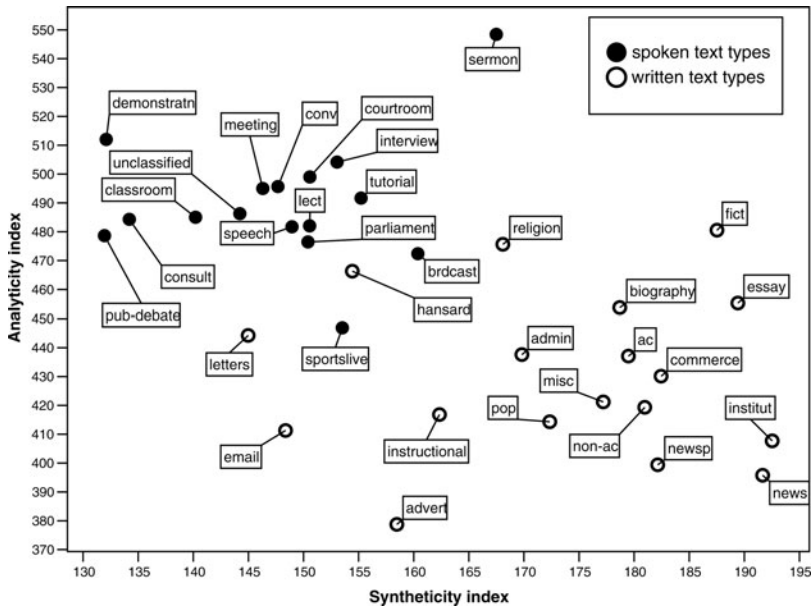


FIGURE 2. BNC macro registers: analyticity by syntheticity (in index points, ptw).

spectrum, we find *public debate* and *demonstrations* as the least synthetic text types in the BNC—genres, therefore, that are least subject to pressures of output economy.

The extreme genres in the analyticity dimension are *sermons* and *advertisements*. Sermons, for one thing, are extremely analytic (AI: 548.4); thus here we are dealing with a text type where the need for explicitness, transparency, and ease of comprehension is rather imperative. (3) exemplifies this genre:

(3) Why not have the light within you so you don't have to go and get it outside but it's there dwelling within you, day by day, moment by moment? And he longs to meet this woman's need. And we can try all sorts of things. And there's, there's things are not necessarily wrong, there's the legitimate things, erm, wi within our work, th there's a, there's job satisfaction, but there's more to that than, in life than just job satisfaction. (BNC text KN8)

What we find in (3), then, is a relatively high degree of reference tracking via pronouns (*you, it, he, we*), many prepositions (e.g., *within, by, in*), and repetition of analytic material galore (notice, for instance, the multiple repetition of existential/dummy *there*). Compare this, now, to (4), an actual advertisement illustrating the BNC's least analytic text type (AI: 378.8):

(4) Build up a total heating system room by room. Interested? USE THE POST-FREE COUPON OVERLEAF. Total Heating. Forget fuel deliveries, dust, dirt,

721 smells, noise, fetching, carrying, tending the boiler . Get a new electric boiler and
 722 forget it—all of it! (BNC text HT1)

723
 724 In (4), it is obvious that all nonessential material is dispensed with, which is, of
 725 course, thanks to the fact that advertisements constitute one of the genres where
 726 the pressure to maximize output economy can be quantified, as it were, in
 727 monetary terms. This pressure affects analytic material in particular, because
 728 analytic markers are typically less compact and economical than synthetic markers.

729
 730
 731

732 *The written-spoken dichotomy*

733 As for higher-order generalizations, let us begin by noting that there are significant
 734 correlations between the AI/SI levels for individual text types and some of the
 735 dimensions of register variation identified by Biber (1988). The relevant
 736 dimensions are *involved vs. informational production* and *abstract vs.*
 737 *nonabstract information*. Based on a subsample of BNC registers¹⁵ whose AIs/
 738 SIs were matched against the factor loadings reported in Biber (1988), the
 739 following pattern emerges. Increased analyticity correlates with involved
 740 production ($r = .78, p < .05$), whereas increased syntheticity dovetails with
 741 abstract informational content ($r = .62, p < .05$). However, it is likely that these
 742 correlations are actually epiphenomenal to a number of very robust differences
 743 between spoken and written text types. These we shall survey in the following text.

744 Some readers will no doubt have noticed already that the z scores in Table 5 are
 745 fairly suggestive with regard to medium: spoken macro registers are typically
 746 associated with positive AI scores and negative SI scores, whereas the converse
 747 holds true for written registers. The box plot in Figure 3 is a more refined way to
 748 look at the variance between spoken and written macro registers (cf. Gries,
 749 2006).¹⁶ In the plot, the boxes depict the interquartile index range comprising
 750 the middle 50% of individual BNC texts (in terms of their analyticity/
 751 syntheticity/grammaticity levels), with the thick line in the boxes indicating the
 752 median. The whiskers above and below the boxes extend to data points that
 753 score no more than 1.5 times the interquartile range. The dots above and below
 754 the whiskers represent outliers, asterisks indicate extreme cases. Four
 755 observations about the variance between spoken and written text types merit
 756 attention:

757
 758
 759
 760
 761
 762
 763
 764
 765

1. *Spoken texts are significantly more analytic than written texts are.* The typical spoken text exhibits 50 or more analytic markers per 1,000 words of running text than the typical written text. In keeping with the interpretational framework outlined earlier, this means that spoken English places a premium on explicitness, transparency, and the minimization of comprehension complexities.
2. *Written texts are significantly more synthetic than spoken texts are,* in that the former exhibit, on average, approximately 30 more synthetic markers per 1,000

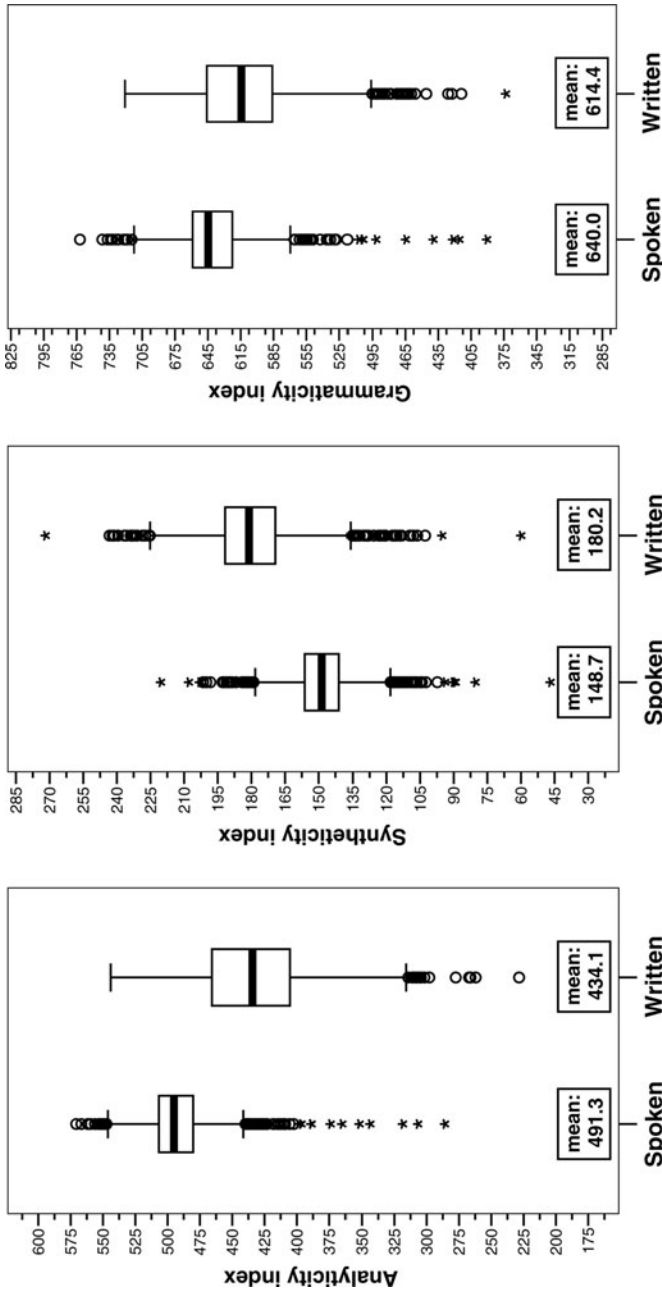


FIGURE 3. Spoken vs. written text types (variance in index points, ptw).

766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810

- 811 words of running text than the latter. Therefore, written texts maximize output
 812 economy whereas spoken texts incur output diseconomies.
- 813 3. *Spoken texts exhibit significantly more grammaticity than written texts.* Thus, vis-
 814 à-vis written texts, spoken texts display more grammatical redundancy, which
 815 eases pragmatic inference complexity.
- 816 4. As far as the scope of variability is concerned, *variability among written texts is*
 817 *more sizable than among spoken texts* (notice the size of the boxes in Figure 3).
 818 For instance, in terms of grammatical analyticity, the interquartile range
 819 containing the middle 50% of all written texts spans roughly 50 index points,
 820 whereas the corresponding interquartile range for spoken texts spans only
 821 about 25 index points.
 822

823 In short, the overall pattern is that spoken texts—which, in Chafe’s (1982) and
 824 Biber’s (1988) parlance, are typically specimens of rather *involved* production—
 825 consistently maximize transparency, explicitness, and ease of comprehension,
 826 whereas written texts, which typically convey more *abstract* (Biber, 1988) or
 827 *detached* (Chafe, 1982) information, can flexibly maximize output economy.
 828 This is due to a crucial and well-known difference between the two mediums:
 829 “Speaking takes place on the fly, but a writer can mull over how best to say
 830 what is desired, and has ample time to edit what is produced”; (Chafe,
 831 1982:262). By the same token, oral comprehension takes place on the fly, with
 832 the spoken word fading rapidly (cf. Hockett, 1960), whereas readers have ample
 833 time to go back and forth in a written text, rereading passages as necessary. The
 834 point is that speech is subject to temporal constraints (specifically transitoriness,
 835 irreversibility, and synchronization, according to Auer, 2009) in a way that
 836 writing is not. This is why, for the sake of comprehension, speakers have
 837 arguably less leeway to manipulate the coding of grammatical information than
 838 writers do. The net result is a more narrow emphasis on transparency and
 839 explicitness in speech, whereas writing “mold[s] a succession of ideas into a
 840 *more complex*, coherent, integrated whole, making use of devices we seldom use
 841 in speaking” (Chafe, 1982:37; emphasis mine).

842 An issue that should also be addressed here is how syntheticity and analyticity
 843 correlate with each other. Earlier, this study detailed that on the level of geographic
 844 varieties, analyticity and syntheticity actually correlate positively, which is contrary
 845 to expectations. Text type variability is similar in this sense. A cursory glance at
 846 Figure 2 might appear to suggest that text types, in fact, exhibit the textbook
 847 trade-off. Globally, increased analyticity incurs reduced syntheticity, and vice
 848 versa. Statistical analysis shows that this relationship has a moderate strength
 849 ($r = -.31$) and is statistically highly significant at $p < .001$. However, it is
 850 important to note that this overall negative correlation disappears entirely—and
 851 thus, turns out to be epiphenomenal—if spoken and written text types are looked
 852 at separately. Among written text types, there is no significant relationship ($r =$
 853 $-.01$, $p = .76$) at all, but among spoken text types there is a weakly *positive*
 854 relationship ($r = .13$, $p < .001$). This replicates this study’s earlier finding on the
 855 geographic plane that was based on spoken data as well.

856 *The sources of text type–stratified analyticity-syntheticity*
 857 *variability*

858 Which grammatical markers or marking families (cf. Tables 1 and 2) are most
 859 involved in the variability just discussed? Starting with analyticity, Table 6 lists
 860 the top five component categories (cf. Tables 1 and 2) that correlate most highly
 861 with overall AI levels.¹⁷ These are (in descending order of importance)
 862 pronouns, as in (5),¹⁸ the negator *not* (contracted or uncontracted), as in (6),
 863 auxiliary *do*, as in (7), modal verbs, as in (8), and auxiliary *have*, as in (9).
 864

- 865 (5) As *she* leaned into the car, the attacker grabbed *her* ... (BNC J1M)
 866 (6) He's *not* out to break any records ... (BNC K1M)
 867 (7) *Does* it work? (BNC K1B)
 868 (8) Shoppers in Abingdon *must* be hoping an agreement is reached ... (BNC K1C)
 869 (9) Eleven people *have* been taken to hospital ... (BNC J1M)
 870

871 Text frequencies of such items, in other words, are the best predictors for overall
 872 analyticity levels. Take, for instance, auxiliary *do*, which has a mean frequency
 873 of 3.7 per thousand words (henceforth: ptw) in the BNC as a whole. In sermons
 874 (the most analytic text type in the corpus), it has a frequency of 8.0 ptw, whereas
 875 in advertisements (the least analytic text type), it has a frequency of merely 1.4 ptw.

876 Turning to syntheticity, Table 7 indicates the top five component categories
 877 correlating most strongly with overall syntheticity levels (again, in descending
 878

879 TABLE 6. *Top five correlations between the analyticity index and broad component*
 880 *categories on the level of individual BNC texts (N = 4,052)*

	Pearson correlation coefficient (<i>r</i>)
883 Pronouns	.75
884 Negator <i>not</i> , <i>n't</i>	.70
885 Auxiliary <i>DO</i>	.64
886 Modals	.51
887 Auxiliary <i>HAVE</i>	.49

888
 889 *Note:* All correlations are significant at $p < .001$.

890 TABLE 7. *Top five correlations between the syntheticity index and broad component POS*
 891 *categories on the level of individual BNC texts (N = 4,052)*

	Pearson correlation coefficient (<i>r</i>)
895 Plural nouns	.59
896 Inflected verbs	.40
897 Conjunctions, subjunctions, prepositions	.39
897 <i>s</i> -genitive	.34
898 Comparative and superlative adjectives	.25

899
 900 *Note:* All correlations are significant at $p < .001$.

901 order of importance): plural nouns, as in (10); inflected verbs, as in (11);
 902 conjunctions, as in (12a), subjunctions, as in (12b), and prepositions, as in (12c);
 903 the *s*-genitive, as in (13); and comparative/superlative adjectives, as in (14).¹⁹

- 904
- 905 (10) Two police armoured *cars* stood outside the courthouse. (BNC A95)
- 906 (11) The US *gave* no answer to their request, *said* Mr Cheney. (BNC A2X)
- 907 (12) a. No record of the initial request is kept *and* the shape and style only evolves
 908 as the metal is worked. (BNC FE6)
- 909 b. Thousands of Soviet television viewers yesterday heard Boris Yeltsin, the
 910 Communist Party rebel, warn of a revolution from below *if* radical
 911 economic changes did not happen within a year (BNC A1G)
- 912
- 913 c. ... the imminent collapse *of* the military regime ... (BNC A1G)
- 914 (13) ... the exiled ANC's internal representatives ... (BNC A1G)
- 915 (14) a. ... the need for *better* resources management. (BNC A96)
- 916
- 917 b. ... the *biggest* burdens in the business ... (BNC A3W)
- 918

919 Notice here, for example, that plural nouns occur with an overall text frequency of
 920 50.8 ptw in the entire BNC—yet in institutional documents (the most synthetic
 921 genre in the BNC), they occur 88.1 times ptw, whereas in public debate (the
 922 least synthetic genre in the BNC), they only have a text frequency of 33.7 ptw. It
 923 is hardly surprising that plural nouns, inflected verbs, the *s*-genitive, and
 924 inflected adjectives are responsible for a considerable amount of syntheticity
 925 variability. What is remarkable is that conjunctions, subjunctions, and
 926 prepositions show up on the list; here we have a per se analytic category, which
 927 nevertheless correlates with increased syntheticity. Why is this? A closer look at
 928 the data reveals that especially the preposition *of* (POS tag PRF; $r = .33$,
 929 $p < .0001$) and other prepositions such as *about*, *at*, *in*, *on*, *on behalf of*, *with*
 930 (POS tag PRP; $r = .46$, $p < .0001$)²⁰ correlate highly with SI levels.²¹ The likely
 931 explanation is that prepositions, although analytic, always come with NPs that
 932 stand a good chance of containing an inflected plural noun, as in (15), or even
 933 premodified by an inflected adjective, as in (16). The net effect is an increase in
 934 syntheticity.

- 935
- 936 (15) ... sell the policy package to voters without worrying *about splits*. (BNC A1J)
- 937 (16) THE arrival in Romania of Mr Gyula Horn, the Hungarian Foreign Minister, is
 938 a sign of Hungarian hopes *for better relations* with their neighbour ... (BNC
 939 AAT)
- 940
- 941
- 942

943 *Interim summary*

944 This section has suggested that there can be a good deal of text type variability
 945 *within* a single geographic variety (in our case, BrE). We have also seen that

946 index levels are predicted by functional pressures and communicative needs. First
 947 and foremost, there is an empirically very robust opposition between spoken and
 948 written texts such that spoken texts exhibit more analyticity as well as
 949 grammaticity, but less syntheticity than written texts. A further difference
 950 between spoken and written English concerns the correlation between analyticity
 951 and syntheticity. Among spoken texts, there is a positive correlation between
 952 analyticity and syntheticity. Among written texts, there is no such correlation.
 953 This section has argued that all these contrasts boil down to the online nature of
 954 speech. Finally, the grammatical categories that cause most of the variability in
 955 the analyticity dimension include pronouns, negators, auxiliary *do/have*, and
 956 modals. The categories that correlate with increased syntheticity comprise—in
 957 addition to the usual suspects (plural nouns, inflected verbs, the *s*-genitive, and
 958 inflected adjectives)—prepositions, which, in spite of being an analytic category
 959 per se, typically attract NPs and the inflectional marking that comes with them.

960

961

SHORT-TERM DIACHRONIC VARIABILITY

962

963

964

965

966

967

968

969

970

971

Adding a longitudinal dimension to the so far purely synchronic discussion, the final parameter of intralingual variability in English to be discussed in this article is real time. More specifically, this section explores short-term diachronic drifts in written English, based on the Brown family of corpora, a set of four matching text corpora documenting early 1960s and early 1990s English, both American and British.

972

973

974

975

976

977

978

979

The Brown family of corpora: an overview

980

981

982

983

984

985

986

987

988

989

990

Table 8 displays global indices in the Brown family of corpora. As for analyticity, notice that there have been significant decreases both in AmE (−8.2 index points; significant at $p = .001^{22}$) and in BrE (−14.8 index points; $p < .001$). The opposite is true for syntheticity. Both matching corpus pairs show significant increases, by 12.7 index points in AmE ($p < .001$) and 9.0 index points in BrE ($p < .001$). We thus note that on the whole, both American and British English have become more synthetic and less analytic over the past half century or so, thus

TABLE 8. *Summary measures for variability in the Brown family of corpora: mean index values and standard deviations (each corpus spans N = 500 texts)*

	Analyticity index		Syntheticity index		Grammaticicity index	
	Mean	SD	Mean	SD	Mean	SD
Brown (1960s AmE)	426.7	37.5	168.9	20.9	595.6	40.9
LOB (1960s BrE)	446.5	33.4	170.7	20.1	617.2	37.0
Frown (1990s AmE)	418.5	42.8	181.6	21.5	600.1	47.7
F-LOB (1990s BrE)	431.6	40.8	179.7	21.8	611.3	46.6

reversing what is often argued to be a millennium-old trend. In terms of grammaticity, no significant changes are observable in AmE, but we note a weakly significant decrease (-5.8 index points; $p = .03$) in BrE. A closer look at the standard deviations provided in Table 8 is also instructive. In both varieties and for all three indices under study, the standard deviations are larger—and sometimes considerably larger—in the 1990s than they were in the 1960s. Thus, written English has come to display more intertextual or inter-register variability in the 1990s than it did in the 1960s. The ensuing discussion will attempt to shed more light on these developments.

Register variability in the Brown family of corpora

We will now scrutinize diachronic drifts in individual written macro and micro registers sampled in the Brown family of corpora. We begin by discussing diachronic drifts among the 15 micro registers sampled in the corpus suite. For every one of these registers, Table 9 displays statistically significant longitudinal AI/SI differentials by national variety (AmE vs. BrE). Let us discuss analyticity and syntheticity variability in turn:

Analyticity. In both BrE and AmE, Press Editorials (category B), Popular Lore (category F), Belles Lettres, Biographies and Essays (category G), the Miscellaneous Learned Writing category (H), and Science (category J) exhibit significant decreases in analyticity. The decreases are most pronounced in categories H and J (Miscellaneous Learned Writing and Science). In AmE, most of the fiction registers—General Fiction (category K), Mystery etc. (category L), Adventure and Western (category N), Romance and Love Story (category P), and Humor (category R)—are on record with increases in analyticity. Among BrE fiction registers, only Romance and Love Story (category P) features a significant positive differential.

Syntheticity. All significant differentials have a positive sign, thus all significant syntheticity differentials are increases. Observe that in both AmE and BrE, the most significant increases have occurred in the Press Reportage section (category A).

As a first step toward a robust generalization, we conflate the 15 micro registers into 4 macro registers. The result of this exercise is shown in Table 10, which displays significant AI and SI differentials for each one of the four macro registers (Press, General Prose, Learned Writing, Fiction). The overall pattern that emerges from the numbers can be summarized as follows. In the analyticity dimension, Press, General Prose, and Learned Writing tend to show significant decreases (which are most substantial in the Learned Writing section), whereas Fiction exhibits significant analyticity increases in both AmE and BrE. As far as syntheticity is concerned, all macro registers except AmE Learned Writing (where the positive differential is not statistically significant) exhibit increases, most markedly so Press language.

The foregoing discussion points to a robust longitudinal split between *informative prose* (registers A through J) and *imaginative prose* (registers K

TABLE 9. Diachronic shifts (in index points, *ptw*) in micro registers in the Brown family of corpora

		Frown vs. Brown (AmE)					F-LOB vs. LOB (BrE)				
		<i>N</i>	AI		SI		<i>N</i>	AI		SI	
A	Press: Reportage	44	+13.0	*	+26.5	***	44			+18.8	***
B	Press: Editorial	27	-11.1	*	+21.9	***	27	-18.6	**		
C	Press: Reviews	17			+14.4	**	17			+9.8	*
D	Religion	17	-22.2	*			17				
E	Skills, Trades, Hobbies	36			+16.1	**	38	-21.0	*		
F	Popular Lore	48	-16.9	*			44	-17.3	*		
G	Belles Lettres, etc.	75	-17.7	***	+12.7	***	77	-22.3	***		
H	Miscellaneous	30	-29.3	**	+13.6	*	30	-38.0	***		
J	Science	80	-36.3	***			80	-32.5	***		
K	General Fiction	29	+19.6	**			29			+11.4	*
L	Mystery etc.	24	+21.7	**	+25.7	***	24				
M	Science Fiction	6			+20.1	*	6				
N	Adventure and Western	29	+14.2	*			29			+13.1	**
P	Romance and Love Story	29	+18.9	*	+16.9	***	29	+13.9	*	+15.7	***
R	Humor	9	+39.2	*			9				

Nonsignificant differentials are omitted. *significant at $p < .05$, **significant at $p < .005$, ***significant at $p < .001$.

1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080

TABLE 10. *Diachronic shifts (in index points, ptw) in macro registers in the Brown family of corpora*

	N	Frown vs. Brown (AmE)				F-LOB vs. LOB (BrE)			
		AI		SI		AI		SI	
Press	88			+22.7	***	-12.2	**	+13.2	***
General Prose	176	-17.6	***	+12.1	***	-20.3	**	+6.3	**
Learned Writing	110	-34.4	**			-34.0	***	+7.2	*
Fiction	126	+20.3	***	+12.7	***	+7.7	*	+11.5	***

Nonsignificant differentials are omitted. *significant at $p < .05$, **significant at $p < .005$, ***significant at $p < .001$.

through R). In Figure 4, we find a diagram that visualizes short-term diachronic drifts among the two text categories in a two-dimensional analyticity-syntheticity plane. The scatterplot makes amply clear that since the 1960s, there has been a pattern of longitudinal divergence between informative and imaginative texts. Both text categories have become more synthetic, but whereas informative prose has become less analytic, imaginative prose has actually become more analytic over time. In other words, written English text types have become more heterogeneous, a fact which partially explains the increasing corpus-internal standard deviations noted earlier in connection with Table 8. The interpretation that the present study would like to offer is that informative prose has traded output economy against explicitness and transparency, thus incurring reader comprehension complexity, whereas imaginative prose has come to favor more grammatical marking, and thus redundancy.

The pattern of divergence between the two text categories is further highlighted when one explores what has happened to grammaticity levels in the period between the early 1960s and the early 1990s. Figure 5 plots the mean difference in grammaticity between sampling times (1990s vs. 1960s) by variety and text category (informative vs. imaginative). We observe that in both varieties, informative prose has shed grammaticity—thus, by inference, eliminating redundancies and maximizing writer output economy at the expense of increased pragmatic complexity whereas imaginative prose has come to be grammatically more redundant (considerably more so).

How can we account for this clear pattern of increasing dissimilarity between informative and imaginative texts? The present study would like to offer that the pattern can be traced back to two tendencies, *economization* and *colloquialization*, which according to the literature are shaping present-day written English. On the one hand, Biber, among others, has noted that modernity has caused an “informational explosion” (Biber, 2003:180), which is why informative texts (e.g., newspaper texts, scholarly prose) are subject to an increasingly high informational density. In this light, Hinrichs & Szmrecsanyi (2007:469) defined *economization* as a tendency toward brevity and compact (grammatical) marking caused by the growing demands of economy and

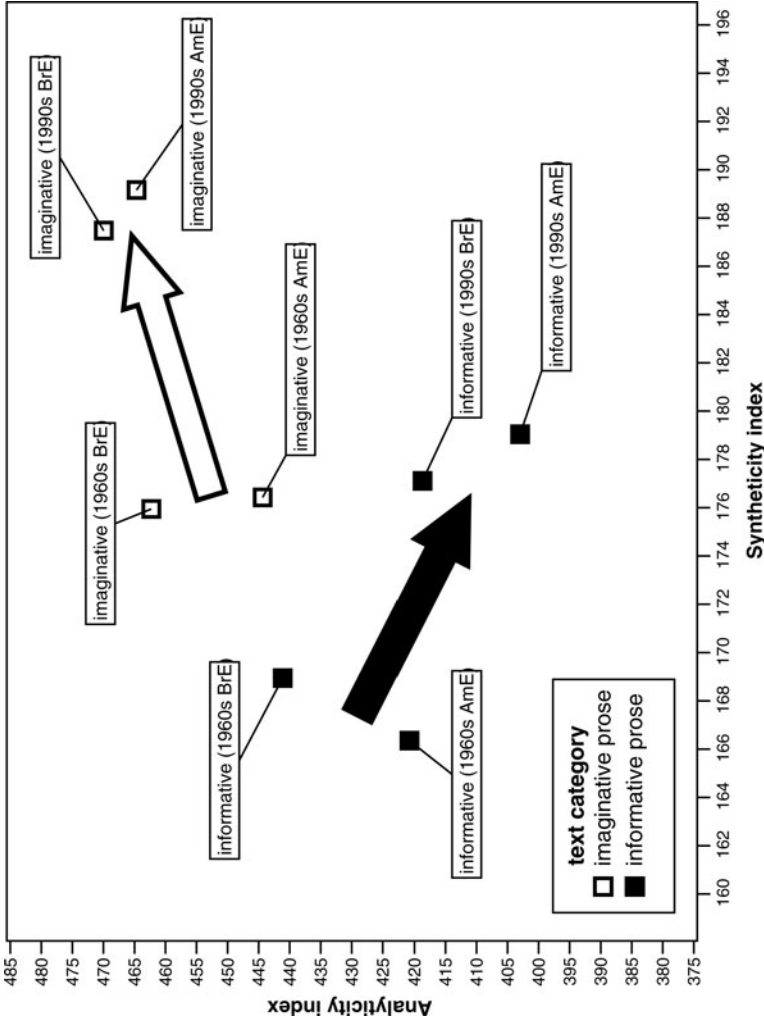


FIGURE 4. Short-term diachronic drifts: analyticity by synthetcity among text categories in the Brown family of corpora (in index points, ptw). All drifts are significant at $p < .05$ in both dimensions.

1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170

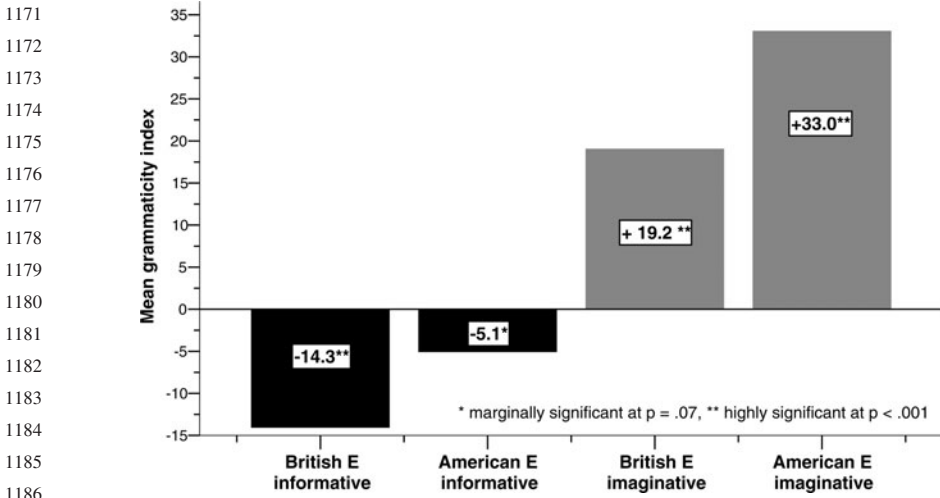


FIGURE 5. Short-term diachronic drifts in the Brown family of corpora: mean increase in grammaticity in text categories (in index points, ptw).

informational compression. On the other hand, it has been observed that since the beginning of the 19th century, certain written genres, such as fiction, have increasingly resembled oral genres (Biber, 2003:169). For the period between the 1960s and the 1990s specifically, there is ample evidence (see, for instance, Hundt & Mair, 1999; Mair, 1997; Mair & Hundt, 1997) that written English has increasingly incorporated more oral features. Against this backdrop, Mair has defined colloquialization as “a trend towards informality ... [which] has had a clear linguistic correlate, a narrowing of the stylistic gap between speech and writing” (Mair, 2006:183). It is important to note, in this connection, that economization and colloquialization are not necessarily conflicting factors. For instance, *not*-contraction is both colloquial and economical. Still, in many cases, the two tendencies are in conflict.

Recall, now, that the section on the written-spoken dichotomy showed that the crucial difference is that spoken texts tend to exhibit more grammaticity and analyticity than written texts do. The increased grammaticity and analyticity of imaginative prose can thus be interpreted as a process of colloquialization such that imaginative written texts narrow the gap to oral texts in two crucial dimensions. Meanwhile, the decreasing grammaticity and analyticity levels observable in informative prose result in better output economy and reduced explicitness and transparency, respectively. This is why the latter drifts practically advertise themselves to be explained in terms of an economization process such that writers seek “to pack information into relatively few words” (Biber, 2003:179).

Having said that, it should be noted that the neat pattern of economization in informative texts and colloquialization in imaginative texts sketched so far does

1216 not fully explain why even imaginative prose has become more synthetic (which
 1217 cannot be a colloquialization phenomenon, because spoken texts are typically *less*
 1218 synthetic than written texts). However, note that we do not really know at this time
 1219 whether spoken language is also becoming more synthetic in real time, an issue
 1220 whose empirical discussion is reserved for another occasion. Pending such
 1221 clarification, we note that colloquialization and economization as processes driving
 1222 change retain a good deal of explanatory potency, despite some interpretational
 1223 twilight concerning what is happening in the syntheticity dimension.

1224
 1225

1226 *The sources of short-term diachronic analyticity-syntheticity*
 1227 *variability*

1228
 1229 Next, we take a more in-depth look at those individual grammatical markers and
 1230 marker families (cf. Tables 1 and 2) that are responsible for the increasing
 1231 divergence between informative and imaginative prose in the Brown family of
 1232 corpora. We start by investigating changes in informative prose. Table 11 lists
 1233 those component categories whose text frequency has significantly²³ changed in
 1234 at least one of the two national varieties under study. An inspection of the signs
 1235 in Table 11 reveals that those categories that have been subject to a frequency
 1236 decrease are typically analytic categories, whereas those categories that show
 1237 increases are typically synthetic in nature. Among the analytic components, it is
 1238 determiners, articles, and *wh*-words that, as a group, show the strongest decrease
 1239 over time. Supplemental analyses suggest that within this category, it is
 1240 especially articles (POS-tag AT), as in (17), and determiners (POS-tag D*), as in
 1241 (18), that show the most substantial frequency decreases over time. Next, with a
 1242 similarly substantial decrease, we find conjunctions, subjunctions, and
 1243 prepositions, a category in which prepositions (POS tag I*), as in (19), are
 1244 responsible for the bulk of the net loss in frequency over time. Also on the
 1245 decline in informative prose is auxiliary *be*, as in (20) and, in the AmE data at
 1246 least, existential *there*, as in (21).

1247

1248 TABLE 11. *Component categories: significant changes (1990s vs. 1960s, in index points ptw)*
 1249 *in informative prose*

	British English (F-LOB vs. LOB)		American English (Frown vs. Brown)		
1253	Determiners, articles, <i>wh</i> -words	-9.1	***	-9.0	***
1254	Conjunctions, subjunctions, prepositions	-7.8	***	-6.0	***
1255	Auxiliary BE			-2.3	***
1256	Existential <i>there</i>			-.5	***
1257	Auxiliary DO			+.4	**
1258	<i>s</i> -genitive	+1.6	***	+3.0	***
1258	Plural nouns	+5.1	**	+7.6	***

1259

1260 *significant at $p < .003$, **significant at $p < .0005$, ***significant at $p < .0001$.

- 1261 (17) In *the* near future, ... (FROWN A12)
- 1262 (18) It won't do him *any* good. (FROWN A01)
- 1263 (19) Polanski now lives *in* Paris, ... (FROWN A24)
- 1264 (20) These fires *were* originally set by lightning or Indians. (FROWN A32)
- 1265 (21) *There* is a common misconception that ... (FROWN A18)

1266
1267 Among those categories that have significantly increased their text frequency in
1268 informative prose over time, notably absent are inflected verbs (cf. Mair, Hundt,
1269 Leech, & Smith, 2002, for a similar finding as to the overall stability of verb
1270 frequencies). Instead, we find the *s*-genitive, as in (22), and especially plural
1271 nouns, as in (23). The prominence of the *s*-genitive as a category on the rise in
1272 informative prose ties in nicely with previous claims (Hinrichs & Szmrecsanyi,
1273 2007; Szmrecsanyi & Hinrichs, 2008) that journalists have come to prefer the
1274 *s*-genitive over the *of*-genitive for primarily output-economy related reasons.
1275 On the other hand, the increasing frequency of plural nouns could conceivably
1276 be seen as supporting evidence for claims that written English has been suffering
1277 from “noun disease” (Potter, 1975:101) of late.

- 1278
- 1279 (22) For Piaget’s constructivist theory ... (F-LOB J23)
- 1280 (23) It is difficult to think productively about ‘modernization’ for many *reasons* ...
- 1281 (F-LOB J26)
- 1282

1283 We now turn to imaginative prose, where the picture is less clear-cut than in the
1284 informative material. Table 12 highlights those component categories whose text
1285 frequency is significantly²⁴ higher in 1990s imaginative prose than it was in
1286 1960s imaginative prose (in imaginative material, no decreases in text frequency
1287 are on record). In AmE, these are, in descending order of importance, pronouns,
1288 as in (24); inflected verbs; as well as determiners, articles, and *wh*-words.
1289 Supplemental analyses show that among inflected verbs, it is the *-s* form of
1290 verbs (POS tag V*Z), as in (25), that is responsible for the overall increase;
1291 among determiners, articles, and *wh*-words, it is possessive determiners (POS
1292 tag APPGE), as in (26), and to a lesser extent what the corpus manual labels as
1293 “*wh*-general adverb[s]” (Hinrichs et al., 2007:24) (POS tag RRQ), as in (27),
1294

1295
1296 TABLE 12. *Component categories: significant changes (1990s vs. 1960s, in index points ptw)*
1297 *in imaginative prose*

	British English (F-LOB vs. LOB)	American English (Frown vs. Brown)	
1301		+11.5	**
1302		+9.5	***
1303		+6.1	*
1304	+5.5		***

1305 *significant at $p < .003$, **significant at $p < .0005$, ***significant at $p < .0001$.

1306 that increase in frequency. In the British imaginative data, finally, it is plural nouns
 1307 (as in (23)) that have been subject to significant expansion.

1308

1309 (24) When *I* reached twenty, *I* moved to New York ... (FROWN R01)

1310 (25) A kid like you *buys* a car ... (FROWN R01)

1311 (26) At that same time *my* father began sending me thick envelopes ... (FROWN
 1312 R06)

1313 (27) ... they walked several blocks to a part of the neighborhood *where* nobody
 1314 knew her. (FROWN R08)

1315

1316

1317

1318

Interim summary

1319 We have seen in this section that written English, both American and British, has
 1320 demonstrably become more synthetic over the past 40 or so years, reversing to some
 1321 extent a millennium-old trend toward more analyticity. In informative prose
 1322 specifically, the component categories that drive the expansion of synthetic
 1323 marking are the inflectional *s*-genitive and inflected plural nouns, but not
 1324 inflected verbs. We have also uncovered evidence that a longitudinal divergence
 1325 between informative and imaginative texts may be taking place, in that there has
 1326 been a development toward less overall analyticity and grammaticity in
 1327 informative prose, whereas imaginative prose shows the converse development.
 1328 The present study has suggested that informative prose is subject to a process
 1329 of economization whereas imaginative prose is undergoing a process of
 1330 colloquialization.

1331

1332

CONCLUSION

1334

1335 The cumulative weight of evidence discussed in this study suggests that English is
 1336 anything but a monolithically analytic, or monolithically nonsynthetic, language.
 1337 Instead, we have seen that observable analyticity, syntheticity, and grammaticity
 1338 levels vary along at least three important dimensions. There is a good deal of
 1339 geographic variation (where sociohistory and variety type seem to impact
 1340 variability), we see significant short-term diachronic variation (where real-time
 1341 variability is induced by changing discourse norms), and the data attest to
 1342 pervasive text type variation (where, among other things, the orality-literacy
 1343 divide plays a major role). In short, we are dealing with variability galore, which
 1344 is demonstrably sensitive to language-external factors. Hence, point estimates
 1345 for, say, “the English language,” which often take center stage in language
 1346 typology, are perhaps more simplistic than is desirable.

1347 On the interpretational plane, this study has linked typological notions to
 1348 language complexity, arguing that grammatical syntheticity and analyticity each
 1349 afford certain payoffs, such as increased explicitness in the case of analyticity
 1350 and better output economy in the case of syntheticity. Against this backdrop,

1351 variability was interpreted in terms of how speakers and writers seek to achieve
 1352 communicative goals while minimizing certain types of complexity (e.g., hearer-
 1353 reader comprehension difficulty) and/or cost (such as the monetary cost
 1354 associated with being exceedingly explicit in advertisements).

1355 Needless to say, there are many more variational dimensions and data sources to
 1356 be investigated in future research. Work is under way in Freiburg to explore long-
 1357 term diachronic analyticity-syntheticity variability in English and to explore
 1358 differences and similarities between non-native, indigenized L2 varieties (such
 1359 as Indian English) and genuine L2 interlanguage varieties (such as French
 1360 learner English). Another dimension of variability that is yet unexplored is how
 1361 sociological variables, such as age, gender, socioeconomic status, might impact
 1362 analyticity-syntheticity variability. Finally, future research along the lines
 1363 sketched out in the present study should also include English-based pidgin and
 1364 creole languages in its data portfolio.

1365 Most pressingly, however, we need data on intralingual variability in other
 1366 languages in order to learn more about the nature of analyticity-syntheticity
 1367 variability and to assess whether the scope of variability observable in English is
 1368 within normal parameters of intralingual variability. Consider text type
 1369 variability. Across all texts in the BNC, the interquartile range (a measure of
 1370 dispersion comprising those 50% of the texts that are closest to the median, thus
 1371 excluding outliers and extreme cases) is 71 index points for the analyticity index
 1372 and 30 index points for the syntheticity index. This is tantamount to saying that
 1373 variability, for example, in analyticity, typically has a scope of ± 35 index points
 1374 in the BNC. The problem is that at present, we have no good idea as to whether
 1375 this range of variability is relatively large in comparison to other languages. Is
 1376 English particularly elastic in regard to analyticity-syntheticity variability? Or
 1377 are there languages (such as Russian, French, Japanese) that are even more
 1378 flexible than English? Is the degree of intralingual elasticity contingent on the
 1379 structural blueprint of the language? An exploration of questions like these
 1380 would admittedly open up an ambitious research agenda, but one that would
 1381 combine careful, intralingual-philological, variationist analysis with the broad,
 1382 abstractive bird's eye perspective that is the hallmark of language typology.
 1383 Indeed, this is an endeavor that would certainly be worth the effort.

1384
 1385 NOTES

- 1386 1. "In Europe, the languages derived from Latin, as well as English, have strongly analytic grammars
 1387 ... synthetic in origin ... they tend strongly toward analytic forms" (translation mine).
 1388 2. The sample size used in Greenberg (1960) were coherent texts of merely 100 words. To mitigate
 1389 the problem of point estimates deriving from such small sample sizes, Stepanov (1995) suggested basing
 1390 the calculation of indices on corpora that "will include hundreds of texts from all existing genres,
 1391 sources, historical periods etc. as one large sample" (Stepanov, 1995:144). Needless to say, this is
 1392 exactly what the present study will do.
 1393 3. The *American National Corpus* is available at: <http://americannationalcorpus.org>.
 1394 4. The rationale behind this choice of dialects is to investigate those counties with the most substantial
 1395 coverage in FRED while maintaining a broad areal coverage. The figures for the subcorpora studied are
 as follows: Somerset—36 interviews, 204,239 words of running text; Kent—11 interviews, 174,420
 words; Shropshire—39 interviews, 174,180 words; Lancashire—23 interviews, 195,111 words;
 Glamorgan—7 interviews, 51,471 words; Sutherland—4 interviews, 10,615 words (Hernández, 2006).

1396 5. The East African data analyzed contain both Kenyan and Tanzanian material, and no distinction
 1397 will be made in what follows between the two varieties.

1398 6. How robust are findings deriving from random samples of 1,000 manually annotated tokens? To
 1399 address this issue, simulations on the basis of the oral history interview material
 1400 (S_interview_oral_history) in the BNC were conducted such that for a number of different sample
 1401 sizes; 10,000 random 1,000-token samples each were obtained to assess the statistical dispersion of
 1402 the mean values for each of the three indices (analyticity, syntheticity, grammaticity) considered in
 1403 the present study. It turns out that a 1,000-tokens random sample has a satisfactorily precise 95%
 1404 confidence interval (CI) for the mean of ± 31 points for the analyticity index, ± 23 points for the
 1405 syntheticity index, and ± 33 points for the grammaticity index. To recapitulate, the former two
 1406 indices span between 0 and 1,000 index points, whereas the latter index spans between 0 and up to
 1407 2,000 points, which means that the 95% CI amounts to less than one one-tenth of 1% of the total
 1408 index range. As for inter-rater reliability of manual annotation based on the extended CLAWS5
 1409 (henceforth CLWAS5e) tag set, parallel annotation by two trained coders of a standard random
 1410 sample data set, drawn from the conversational section of ICE-NZ and spanning $N = 1,000$ tokens,
 1411 yielded a simple agreement rate of approximately 91% and an “excellent” (Orwin, 1994:152)
 1412 Cohen’s κ value of .90.

1409 7. Although the Brown family’s CLAWS8 tag set has special tags for auxiliary usage of primary
 1410 verbs, for the sake of comparability to the BNC, auxiliary usage was identified contextually.

1411 8. In regard to modal verbs (*may—might, will—would*, etc.) and pronouns (*I, you, we*, etc.), note that
 1412 these elements are classified here as categorically analytic tokens, although some analysts would view
 1413 forms such as *might* and *would* and even pronominal forms as inflected elements. The primary reason for
 1414 not letting these elements load on the syntheticity index is that the postulated derivation of, say, the form
 1415 *would* from *will* or the derivation of, say, *we* from *I* is not likely to have the same status—semantically, in
 1416 terms of productivity, and cognitively (on the part of language users)—as the derivation of, say, *sang*
 1417 from *sing* or *houses* from *house*. More pragmatically speaking, varieties of English also do not seem
 1418 to exhibit much variability in regard to the frequency of such tokens. Be that as it may, we concede
 1419 that the exclusion of such tokens from the syntheticity index is ultimately arbitrary; a detailed
 1420 discussion of how an inclusion of these elements may or may not change results is reserved for
 1421 another occasion.

1418 9. This example calculation will illustrate. Assume a text spanning 2,000 running words exhibits 300
 1419 synthetic markers and 800 analytic markers. The resulting indices are calculated as follows: SI: $300/$
 1420 $2000 \times 1000 = 150$; AI: $800/2000 \times 1000 = 400$; GI: $300 + 800/2000 \times 1000 = 550$.

1421 10. Methodologically, this section roughly follows a set of pilot studies (Kortmann & Szmrecsanyi,
 1422 forthcoming; Szmrecsanyi & Kortmann, 2009) on pertinent variability in World Englishes. Note,
 1423 however, that the discussion in the present study is based on a different data set and on a more
 1424 sophisticated coding method.

1424 11. Here and in the following, p values derive from independent samples t tests, unless stated otherwise.
 1425 In all cases, the data are approximately normally distributed (< 2 standard errors of skewness).

1425 12. In this regard, it might be noted that there are moderately strong albeit nonsignificant correlations
 1426 between population size and the indices in question. Population size correlates positively with analyticity
 1427 ($r = .19, p = .49$) and negatively with syntheticity ($r = -.34, p = .19$). The generalization is that large
 1428 speaker communities (where we find more contact) tend to prefer analytic marking, which is more
 1429 explicit and transparent, whereas small speaker communities, where we typically find less dialect
 1430 contact, tend to prefer synthetic marking, which optimizes output economy.

1430 As for population size, the figures (in million inhabitants) entered into analysis are as follows:
 1431 Glamorgan: 2.1; East African E: 1.5; New Zealand E: 3.9; Indian E: 1000.0; Kent: 1.3; Lancashire:
 1432 1.1; Shropshire: .3; Somerset: .5; Sutherland: .2; Irish E: 4.3; Hong Kong E: 6.8; Philippine E: .9;
 1433 Singapore E: 4.2; Jamaican E: 2.6; Standard BrE: 60.2; Standard AmE: 300.0 (source:
 1434 Encyclopaedia Britannica ultimate reference suite DVD. London: Encyclopaedia Britannica, 2004).

1433 13. Given that 15 broad component categories were tested, the Bonferroni-corrected α level calculates
 1434 as $p = .05/15 = .003$.

1435 14. In the same vein, population size, which is arguably a proxy for language/dialect contact, tends to
 1436 correlate positively with analyticity and negatively with syntheticity.

1436 15. More specifically, the following 20 registers were investigated: W_ac_medicine,
 1437 W_ac_tech_engin, W_ac_soc_science, W_ac_nat_science, W_ac_humanities_arts,
 1438 W_newsp_other_report, W_newsp_brdsht_nat_report, W_newsp_brdsht_nat_editorial, W_pop_lore,
 1439 W_letters_prof, W_letters_personal, W_biography, W_religion, W_fiction, S_broadcast,
 1440 S_sportslive, S_speech_scripted, S_speech_unscripted, S_interview, S_conv.

- 1441 **16.** Note that the box plot is based on individual BNC texts. Spoken-written differences in mean index
 1442 values are, according to independent samples *t* tests, highly significant at $p < .001$ throughout. Notice
 1443 also that with skewness values of $< .9$ standard errors of skewness in either dimension, the data are
 1444 approximately normally distributed.
- 1445 **17.** Given that 15 component categories were tested, note that the correlations in Table 6 are significant
 1446 at a Bonferroni-corrected α level of $p = .05/15 = .003$.
- 1447 **18.** We wish to point out in this connection that the variational alternative to pronoun usage is not
 1448 necessarily a full NP (e.g., *Mary leaned into the car* instead of *she leaned into the car*), but possibly
 1449—in certain contexts and registers (for example, conversation)—a null form (e.g., *He can't stand*
 1450 *his mother. Ø Can't say I blame him.* [BNC AC3]). Notice here that a propensity for null subjects
 1451 has also been reported for some regional varieties of English, such as Newfoundland English
 1452 (Wagner 2007).
- 1453 **19.** The correlations in Table 7 are significant at a Bonferroni-corrected α level of $p = .05/15 = .003$.
- 1454 **20.** Given that 55 individual POS tags were tested, the significance levels reported here are significant
 1455 at a Bonferroni-corrected α level of $p = .05/55 = .0009$.
- 1456 **21.** The subordinating conjunction *that* (POS tag: CJT) does not, in fact, correlate significantly with SI;
 1457 coordinating conjunctions (POS tag: CJC) correlate positively, if weakly, with SI ($r = .04, p = .008$),
 1458 whereas subordinating conjunctions other than *that* (POS tag: CJS) actually correlate negatively with
 1459 SI ($r = -.25, p < .0001$).
- 1460 **22.** Here and in the following, p values derive from independent samples *t* tests, which have been run
 1461 on the basis of individual corpus texts ($N = 500$ texts for each corpus). In all cases, the data are
 1462 approximately normally distributed (< 2 standard errors of skewness).
- 1463 **23.** Given that 15 broad component categories were tested, the Bonferroni-corrected α level calculates
 1464 as $p = .05/15 = .003$.
- 1465 **24.** Given that 15 broad component categories were tested, the Bonferroni-corrected α level calculates
 1466 as $p = .05/15 = .003$.

REFERENCES

1463 Anttila, Raimo. (1989). *Historical and Comparative Linguistics*. Philadelphia: Benjamins.

1464 Aston, Guy, & Burnard, Lou. (1998). *The BNC handbook: Exploring the British National Corpus with*
 1465 *SARA*. Edinburgh: Edinburgh University Press.

1466 Auer, Peter. (2009). On-line syntax: Thoughts on the temporality of spoken language. *Language*
 1467 *Sciences* 31(1):1–13.

1468 Biber, Douglas. (1988). *Variation across speech and writing*. Cambridge: Cambridge University Press.

1469 ———. (2003). Compressed noun-phrase structure in newspaper discourse: The competing demands of
 1470 popularization vs. economy. In J. Aitchison & D. M. Lewis (eds.), *New media language*. New York:
 1471 Longman. 169–181.

1472 Bisang, Walter. (2009). On the evolution of complexity—Sometimes less is more in East and mainland
 1473 Southeast Asia. In G. Sampson, D. Gil, & P. Trudgill (eds.), *Language complexity as a variable*
 1474 *concept*. Oxford: Oxford University Press.

1475 Bussmann, Hadumod, Trauth, Gregory, & Kazzazi, Kerstin. (1996). *Routledge dictionary of language*
 1476 *and linguistics*. New York: Routledge.

1477 Chafe, Wallace L. (1982). Integration and involvement in speaking, writing, and oral literature. In
 1478 D. Tannen (ed.), *Spoken and written language: Exploring orality and literacy*. Norwood,
 1479 NJ: Ablex. 35–53.

1480 Danchev, Andrei. (1992). The evidence for analytic and synthetic developments in English. In
 1481 M. Rissanen, O. Ihalainen, T. Nevalainen, & I. Taavitsainen (eds.), *History of Englishes: New*
 1482 *methods and interpretations in historical linguistics*. New York: Mouton de Gruyter. 25–41.

1483 Francis, Nelson W., & Kučera, Henry. (1982). *Frequency analysis of English usage: Lexicon and*
 1484 *grammar*. Boston: Houghton Mifflin.

1485 Greenbaum, Sidney (ed.). (1996). *Comparing English worldwide: The international corpus of English*.
 Oxford: Clarendon.

Greenberg, Joseph H. (1960). A quantitative approach to the morphological typology of language.
International Journal of American Linguistics 26(3):178–194.

Gries, Stefan Th. (2006). Exploring variability within and between corpora: Some methodological
 considerations. *Corpora* 1(2):109–151.

Hernández, Nuria. (2006). *User's guide to FRED*. Freiburg: English Dialects Research Group. Available
 online at: <http://www.freidok.uni-freiburg.de/volltexte/2489/>.

- 1486 Hinrichs, Lars, & Szmrecsanyi, Benedikt. (2007). Recent changes in the function and frequency of
 1487 standard English genitive constructions: A multivariate analysis of tagged corpora. *English*
Language and Linguistics 11(3):437–474.
- 1488 Hinrichs, Lars, Waibel, Birgit, & Smith, Nicholas. (2007). *The POS-tagged, postedited F-LOB and*
 1489 *Frown corpora: A manual, including pointers for successful use*. Freiburg: Department of English,
 1490 University of Freiburg. Available online at: [https://webpace.utexas.edu/lh9896/public/hinrichs/](https://webpace.utexas.edu/lh9896/public/hinrichs/Manual_final.pdf)
 1491 [Manual_final.pdf](https://webpace.utexas.edu/lh9896/public/hinrichs/Manual_final.pdf).
- 1491 Hockett, Charles F. (1954). Two models of grammatical description. *Word* 10:210–231.
- 1492 ———. (1960). The origin of speech. *Scientific American* 203:88–96.
- 1493 Humboldt, Wilhelm von. (1836). *Über die Verschiedenheit des menschlichen Sprachbaues und ihren*
 1494 *Einfluss auf die geistige Entwicklung des Menschengeschlechts*. Berlin: Dümmler.
- 1494 Hundt, Marianne, & Mair, Christian. (1999). ‘Agile’ and ‘uptight’ genres: The corpus-based approach to
 1495 language change in progress. *International Journal of Corpus Linguistics* 4:221–242.
- 1496 Johansson, Stig, & Hofland, Knut. (1989). *Frequency analysis of English vocabulary and grammar.*
 1497 *Based on the LOB Corpus*. Oxford: Clarendon.
- 1497 Kasevič, Vadim, & Jachontov, Sergej E. (eds.) (1982). *Kvantitativnaja tipologija jazykov Azii i Afriki* [A
 1498 quantitative typology of Asian and African Languages]. Leningrad: Izdatel'stvo Leningradskogo
 1499 *universiteta*.
- 1499 Kelemen, József. (1970). Sprachtypologie und Sprachstatistik. In L. Dezso & P. Hajdú (eds.), *Theoretical*
 1500 *problems of typology and the Northern Eurasian languages*. Amsterdam: Gruener. 53–63.
- 1501 Kempgen, Sebastian, & Lehfeldt, Werner. (2004). Quantitative typologie. In G. Booij, C. Lehmann,
 1502 J. Mugdan, & S. Skopeteas (eds.), *Morphologie. Ein internationales Handbuch zur Flexion und*
Wortbildung. Berlin: Mouton de Gruyter. 1235–1246.
- 1503 Kortmann, Bernd, & Szmrecsanyi, Benedikt. (forthcoming). World Englishes between simplification
 1504 and complexification. In L. Siebers & T. Hoffmann (eds.), *World Englishes: Problems—Properties*
 1505 *—Prospects*. Philadelphia: Benjamins.
- 1505 Labov, William. (1966). *The social stratification of English in New York City*. Washington, DC: Center
 1506 for Applied Linguistics.
- 1507 ———. (1972). *Sociolinguistic patterns*. Philadelphia: University of Philadelphia Press.
- 1507 Mair, Christian. (1997). Parallel corpora: A real-time approach to the study of language change in
 1508 progress. In M. Ljung (ed.), *Corpus-based studies in English*. Amsterdam: Rodopi. 195–209.
- 1509 ———. (2006). *Twentieth-century English: History, variation, and standardization*. Cambridge:
 1510 Cambridge University Press.
- 1510 Mair, Christian, & Hundt, Marianne. (1997). The corpus-based approach to language change in
 1511 progress. In H. Sauer & U. Böker (eds.), *Anglistentag 1996: Proceedings*. Tübingen: Niemeyer.
 1512 71–82.
- 1513 Mair, Christian, Hundt, Marianne, Leech, Geoffrey, & Smith, Nicolas. (2002). Short-term diachronic
 1514 shifts in part-of-speech frequencies: A comparison of the tagged LOB and F-LOB. *International*
Journal of Corpus Linguistics 7:245–264.
- 1515 Marty, Anton. (1908). *Untersuchungen zur Grundlegung der allgemeinen Grammatik und*
 1516 *Sprachphilosophie*. Halle a.S.: Niemeyer.
- 1516 Meshrie, Rajend. (2006a). World Englishes and the multilingual history of English. *World Englishes* 25
 1517 (3/4):381–390.
- 1518 ———. (2006b). Anti-deletions in an L2 grammar: A study of Black South African English mesolect.
 1519 *English World-Wide* 27(2):111–145.
- 1520 Orwin, Robert. (1994). Evaluating coding decisions. In H. Cooper & L. Hedges (eds.), *The handbook of*
research synthesis. New York: Russell Sage Foundation. 139–162.
- 1521 Potter, Simeon. (1975). *Changing English*. London: Deutsch.
- 1522 Schlegel, August Wilhelm von. (1818). *Observations sur la langue et la littérature provençales*. Paris:
 1523 Librairie grecque-latine-allemande.
- 1523 ———. (1846). *Œuvres de M. Auguste-Guillaume de Schlegel: Écrites en français et publiées par*
 1524 *Édouard Böcking*. Leipzig: Weidmann.
- 1525 Schwegler, Armin. (1990). *Analyticity and syntheticity: A diachronic perspective with special reference*
to Romance languages. New York: Mouton de Gruyter.
- 1526 Stepanov, Arthur V. (1995). Automatic typological analysis of Semitic morphology. *Journal of*
 1527 *Quantitative Linguistics* 2(2):141–150.
- 1528 Szmrecsanyi, Benedikt, & Hernández, Nuria. (2007). *Manual of information to accompany the Freiburg*
 1529 *Corpus of English Dialects Sampler (“FRED-S”)*. Freiburg: English Dialects Research Group.
 1530 Available online at: <http://www.freidok.uni-freiburg.de/volltexte/2859/>.

1531 Szmrecsanyi, Benedikt, & Hinrichs, Lars. (2008). Probabilistic determinants of genitive variation in
 1532 spoken and written English: A multivariate comparison across time, space, and genres. In
 1533 T. Nevalainen, I. Taavitsainen, P. Pahta, & M. Korhonen (eds.), *The dynamics of linguistic
 1534 variation: Corpus evidence on English past and present*. Amsterdam: Benjamins. 291–309.

1535 Szmrecsanyi, Benedikt, & Kortmann, Bernd. (2009). Between simplification and complexification:
 1536 Non-standard varieties of English around the world. In G. Sampson, D. Gil, & P. Trudgill (eds.),
 1537 *Language complexity as a variable concept*. Oxford: Oxford University Press. 65–79.

1538 Trudgill, Peter. (2001). Contact and simplification: Historical baggage and directionality in linguistic
 1539 change. *Linguistic Typology* 5(2/3):371–374.

1540 ———. (2009). Vernacular universals and the sociolinguistic typology of English dialects.
 1541 In M. Filppula, J. Klemola, & H. Paulasto (eds.), *Vernacular universals and language contacts:
 1542 Evidence from varieties of English and beyond*. London: Routledge. 304–322.

1543 Vennemann, Theo. (1982). Isolation—Agglutination—Flexion? Zur Stimmigkeit typologischer
 1544 parameter. Fakten und theorien. In S. Heinz & U. Wandruszka (eds.), *Festschrift für Helmut Sinn
 1545 zum 65. Geburtstag*. Tübingen: Narr. 327–334.

1546 Wagner, Susanne. (2007). Null subjects in English—economically motivated? Paper presented at the
 1547 36th Conference on New Ways of Analyzing Variation (NWA36). Philadelphia.

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1563

1564

1565

1566

1567

1568

1569

1570

1571

1572

1573

1574

1575