

1 **Predicting Voice Alternation Across Academic Englishes**

2

3 **1. Introduction**

4 One of the aims of the *International Corpus of English (ICE)*¹ is to enable comparative
5 studies of variation across a broad range of World Englishes (WEs). Initially, research
6 was somewhat limited by the fact that most ICE corpora lacked grammatical
7 annotation. This is particularly problematic for the study of relatively abstract
8 grammatical patterns, such as voice. Even in a part-of-speech-tagged corpus, it
9 remains difficult to efficiently retrieve active transitive counterparts of passives,
10 making principled study of the active-passive alternation across WEs practically
11 impossible. Now that a large number of ICE components have been parsed, automatic
12 retrieval of actives and passives is possible.

13 Voice alternation is of interest because previous research (e.g. Leech et al.
14 2009) identified regional differences in a shift from passive towards active
15 constructions, which is more pronounced in American (AmE) than in British English
16 (BrE) academic writing. The question was whether the trend towards greater use of
17 the active voice could also be found in other first- and second-language varieties of
18 English (see Hundt et al. 2016). In this paper, we add to previous research by looking
19 into the factors predicting the choice between active and passive. In particular, we
20 aim to discover whether or not WEs share the core grammar of voice alternation.
21 Substrate influence in academic writing – a highly edited register – will not be the
22 same as in spontaneous spoken language where we find evidence of obvious
23 structural transfer or borrowing, e.g. in the form of the *kena* passive in Singapore
24 English (Bao & Wee 1999). All the same, substrate may play a role at a more subtle
25 level, e.g. in a variety’s readiness to combine passive voice with marking for other
26 verbal categories such as aspect or giving preference to inanimate subjects (see
27 section 2.4).

28 In part two of the paper, we provide background on previous research, the

¹ ICE samples acrolectal/standard English spoken as a first or institutionalised second language. For more background information, see the project website at [URL will be provided once the new website has been announced].

1 rationale for the choice of varieties we investigate, and on passives in the substrate
2 languages. Section three gives information on data retrieval and coding of factors.
3 Section four reports the results of our statistical analysis, which are discussed in the
4 context of modelling probabilistic variation across WEs (section five).

5 **2. Background**

6 **2.1 Previous research: *be*-passives across time and space**

7 The overall frequency of *be*-passives varies across time and space: the long-term trend
8 towards declining use of passives is more advanced in AmE than in BrE (Biber &
9 Finegan 1989; Seoane 2006; Leech et al. 2009). On the basis of manual analysis, earlier
10 studies show that this decline happens at the expense of active transitives (Seoane &
11 Loureiro-Porto 2005) in academic writing, and that the trend is more pronounced in
12 natural sciences than in the humanities (Hundt & Mair 1999). Research on the passive
13 in second-language varieties of English is scarce (but see, e.g., Biewer 2009).

14 Hundt et al.'s (2016) analysis of the academic section of the parsed Brown-
15 family of corpora and 15 ICE components corroborates previous results. Somewhat
16 surprisingly, however, the broader selection of WEs did not reveal a divide into first-
17 (ENL) and second-language (ESL) varieties. Instead, AmE was the only variety that
18 significantly differed from other WEs in preferring actives in academic writing.
19 Regression analysis showed that variation across sub-disciplines (humanities and
20 social sciences vs. natural sciences and technology) was even more pronounced than
21 regional variation, with the 'soft sciences' preferring a more active style (see also Biber
22 & Finegan 1989). In other words, despite substantial differences in the substrate
23 languages, stylistic preferences in the various sub-disciplines turned out to be the
24 decisive factor for the overall frequencies of *be*-passives and active transitives from
25 15 academic Englishes in Africa, America, Asia, Europe and the Pacific. AUTHORS
26 (forthcoming) zoom in on the role that authorial presence plays in the choice of active
27 over passive in six ENL varieties. The focus in this study is on language-internal factors
28 predicting voice alternation.

29 **2.2 Factors predicting choice of passive voice**

30 Several factors are known to determine use of passives in general. In this study, we

1 look at four that lend themselves to modelling in a multivariate analysis, two syntactic,
2 one semantic and one discourse-pragmatic factor.

3 *Complexity of the verb phrase*

4 English is a language that marks certain categories in the VP (voice, aspect, and
5 occasionally also tense) periphrastically. This can give rise to quite complex VPs, such
6 as *they will have been being chased*. While such maximally complex VPs are rare (see
7 Hundt 2013: 167), combinations of the perfect or progressive with passive are used
8 fairly regularly. One of the strategies in second language acquisition (especially by
9 adults) is to reduce structural complexity; but errors may also lead to more complex
10 patterns (see Kortmann & Szmrecsanyi 2009 or Thomason 2013). We would therefore
11 expect to see differences between ENL and ESL varieties in the combination of passive
12 voice with other verbal categories, and a tendency for learners to avoid passives in
13 complex VPs.

14 *Weight*

15 From a syntactic perspective, passives are often seen to rearrange the order of
16 elements so that long, heavy constituents occupy final position, thus conforming to
17 the Principle of End Weight, as first put forward by Behagel (1909, 1930). Since this is
18 likely to be a universal, cognitive processing constraint, we expect constituent weight
19 to play a significant role in the choice between active and passive across ENL and ESL
20 varieties. If anything, the factor might be more marked in the ESL varieties than in BrE
21 and AmE.

22 *Animacy*

23 From a semantic standpoint, English, like many other languages, tends to have human
24 subjects and topics, as encapsulated in the Animacy Hierarchy originally proposed by
25 Silverstein (1976): human > non-human animate > inanimate. Passives have been
26 shown to contravene this hierarchy (Seoane 2009: 375-379). At the same time,
27 animacy as a predictor for syntactic variation has turned out to be subject to regional
28 variation in WEs in various studies that model frequency effects in grammar (e.g.
29 Hinrichs & Szmrecsanyi 2007, Bresnan & Hay 2008 or Bresnan & Ford 2010). It is
30 possible that it affects voice alternation to different degrees across WEs, and that ESL
31 varieties might prefer animate subjects in both actives and passives.

1 *Givenness of the active object and passive subject*

2 In a language with fixed word-order, like English, passives work as a possible order
3 rearrangement strategy to conform to the given-before-new principle, which
4 concerns the degree of accessibility of information to discourse participants. This
5 pragmatic information status (given vs. new) is a scalar concept rather than a
6 dichotomy, and several taxonomies have been proposed to measure the degree of
7 givenness (see Seoane 2012). The accessibility of information depends not only on the
8 linguistic context (what has already been mentioned and therefore activated in the
9 interlocutors' minds) but also on the extra-linguistic context (knowledge of the world,
10 the perception of the immediate extra-linguistic context by interlocutors). Though
11 there is not a one-to-one correlation between givenness and definiteness (definite *the*
12 *chair* is supposed to convey given meaning vs. *a chair* which is typically used to
13 introduce a new referent in discourse) analyzing the degree of definiteness of active
14 object and passive subject can help identify contrasts between given and new
15 elements.

16 **2.3 Englishes selected**

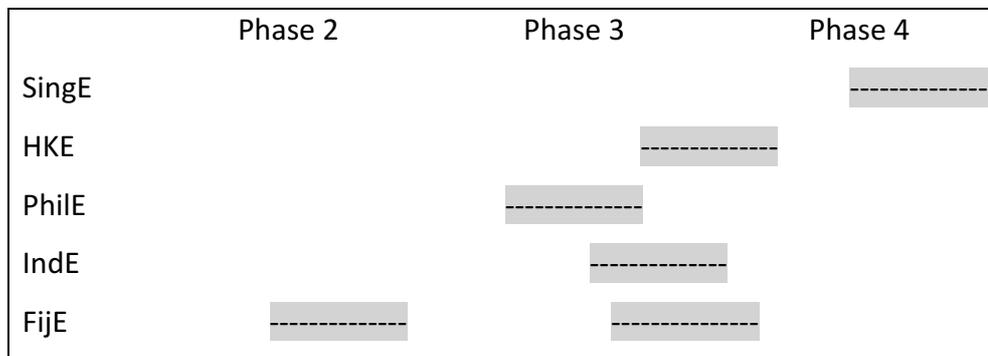
17 This paper aims to add to previous studies by looking at contextual factors which may
18 play a role in voice alternation. We also compare ENL and ESL varieties, thus adding
19 to a growing body of probabilistic research into WEs. The typical pattern to emerge
20 from such studies is that of a shared core probabilistic grammar with variety-specific
21 peculiarities at a more fine-grained level of analysis. Szmrecsanyi et al. (2016: 133) call
22 this 'probabilistic indigenization', which they define as

23 the process whereby stochastic patterns of internal linguistic variation are
24 reshaped by shifting usage frequencies in speakers of post-colonial varieties. To the
25 extent that patterns of variation in a new variety A, e.g. the probability of item x in
26 context y, can be shown to differ from those of the mother variety, we can say that
27 the new pattern represents a novel, if gradient, development in the grammar of A.
28 These patterns need not be consistent or stable ..., but they nonetheless reflect the
29 emergence of a unique, region-specific grammar.

30 For the present study, we selected BrE and AmE as our ENL reference varieties. In our
31 choice of post-colonial contact varieties, we focus on the Asia-Pacific region, selecting
32 Singapore (SingE), Hong Kong (HKE), Indian (IndE) and Fiji English (FijE) as varieties that
33 are historically related to BrE, and Philippine English (Phile) as a second-language

1 variety deriving from AmE.

2 In his seminal work on post-colonial Englishes, Schneider (2007) outlines five
3 developmental stages for the evolution of new Englishes. Figure 1 shows that the
4 second-language varieties we selected for this case study can be found at stages
5 two/three and four along this cycle:



6 **Figure 1: Developmental phase of post-colonial varieties of English according to Schneider's (2007) model**

7 According to Schneider (2007: 114-18), FijE is the least advanced variety here,
8 i.e. it is still at the stage of exo-normative stabilization in a bilingual setting with the
9 original 'homeland' as the norm-giving centre but only incipient structural borrowing;
10 other studies (Geraghty et al. 2006, Zipp 2014, Hundt et al. 2015) indicate that it has
11 progressed to stage three (nativization) and might even be showing the first signs of
12 moving to stage four (endo-normative stabilization). HKE, IndE and PhilE have
13 progressed to the stage of structural nativization, which sees the development of
14 distinctly local structural characteristics; HKE and IndE are moving towards stage four,
15 whereas PhilE is likely to remain at stage three; SingE, finally, has moved on to the
16 stage of endo-normative stabilization, i.e. the emergence of a relatively homogenous
17 local variety that finds acceptance in the speech community (see Schneider 2007).

18 In addition to differences in the development towards new Englishes in
19 Schneider's model, the five contact varieties were selected because the (dominant)
20 substrate languages potentially allow for different outcomes with respect to structural
21 factors that influence alternation between active and passive.

22 **2.4 Passives and language contact**

23 We use the term 'passive' here rather than 'voice' because passive can be grammatical
24 but also semantic and/or pragmatic. In English, for instance, passive voice is marked

1 grammatically with catenatives *be* or *get* followed by a past participle.² English, too,
2 has constructions which can be argued to be notionally passive but which are formally
3 active in terms of ‘voice’: examples would be ‘middles’ (e.g. *The book reads well*) or
4 remnants of the earlier passival (e.g. *Coffee is now serving*) (see Hundt 2004, 2007).

5 According to Keenan (1985: 247), a basic passive is characterised by the
6 following properties: “(i) no agent phrase ... is present, (ii) the main verb (in its non-
7 passive form) is transitive, and (iii) the main verb expresses an activity, taking agent
8 subjects and patient objects.” In other words, basic passives, from a typological point
9 of view, do not mark voice morphologically on the verb. In the following, we look at
10 the main substrate languages for the ESL varieties we investigate to see whether they
11 mark passives morphologically, in other ways, or not at all.

12 *Substrate languages for Singapore and Hong Kong English*

13 The most important substrate languages in these two East Asian countries are
14 Mandarin and Cantonese Chinese, respectively. Other substrate languages relevant
15 for SingE are Malay and Hokkien (and to a lesser extent Tamil, spoken by the Indian
16 population).

17 In Hokkien (a group of Min Nan Chinese dialects), Mandarin (an isolating
18 language of the Sino-Tibetan group) and Malay (an agglutinative language belonging
19 to the Austronesian family of languages) verbs do not inflect for voice. In Mandarin,
20 for example, *bei* has grammaticalised as a “function word with no inherent meaning
21 other than passiveness marking” (Xiao et al. 2006: 125); it can appear in both long and
22 short passives (i.e. those with or without an explicit agent) without any change in the
23 verb form. Other less frequent passive markers in Mandarin are *gei*, *jiao* and *rang*
24 (Xiao et al. 2006: 125). As for the frequency and distribution of the passive in
25 Mandarin, Xiao et al. (2006: 124-141) observe that passives are less frequent in
26 Mandarin than in English because they normally convey negative and adversative
27 meanings; only recently, because of Western influence, have passives started to be
28 used with positive or neutral meanings (Xiao et al. 2006: 135, Gunn 2017). Passives
29 are especially infrequent in academic writing since, unlike in English, they do not mark
30 objectivity and impersonality.

² Only *be* but not *get* also satisfies the criteria for auxiliarihood.

1 The passive in Cantonese is similar to that of Mandarin, the main difference
2 being that the passive marker *bei* can only appear in long passives in Cantonese
3 (Matthews and Yip 1994: 149), i.e. passives with an overtly expressed agent. As is the
4 case in Mandarin, passives in Cantonese are less frequent than in English, because
5 there are other strategies in the language which readily topicalise objects and because
6 they are still strongly associated with the expression of adversative meaning
7 (Matthews and Yip 1994: 150).

8 *Substrate languages for Philippine English*

9 The linguistic situation of the Philippines is one of intense multilingualism (McFarland
10 2008: 143). The most widely spoken indigenous language is Tagalog (and the standard
11 version Filipino since 1987), which – like Malay – is an agglutinative language of the
12 Austronesian family. The voice system in Tagalog is a very controversial topic,
13 especially as far as the status of the grammatical subject and the typological nature of
14 the language are concerned (see Shibatani 1988: 85-142). In brief, Tagalog has a goal-
15 topic construction, in which the verb is marked inflectionally as requiring the
16 patient/goal of the action (and not the agent) to be the topic. In addition, it sometimes
17 takes an actor complement, the semantic equivalent of an English *by*-phrase
18 (Schachter & Otnes 1972: 73; Shibatani 1988: 86-89).

19 *Dominant substrate influence on Indian English: Hindi*

20 In India, too, we are looking at a multilingual situation with various local languages
21 that are likely to have an impact on the English spoken. At the same time, most studies
22 considering substrate influence in IndE typically focus on Hindi, since it is an important
23 indigenous national language. Hindi is a morphologically polysynthetic language of the
24 Indo-Aryan language family. In Hindi, passive VPs consist of the passive auxiliary *ja*
25 following an inflected past participle of the main verb. Most of the passives in Hindi
26 are agentless, their main function being that of eliding an irrelevant agent (Sandahl
27 2000: 101; Kachru 2006: 93).

28 *Substrate and Fiji English*

29 In the Fiji Islands, dialects of Fijian and a non-standard variety of Hindi (Fiji Hindi)
30 comprise the main substrate influence for FijE. We only comment on Fijian here, as
31 passive marking in the local variety of Hindi is similar to what we find in standard Hindi.

1 Fijian is an Austronesian (Malayo-Polynesian) language that uses some inflectional
 2 endings on the verb and is thus somewhat less isolating than other languages in the
 3 South Pacific (Lynch 1998: 130f.). Interestingly, transitivity is one of the categories that
 4 are marked with the help of a suffix Fijian verbs (ibid.: 140). There are three broad
 5 types of transitives in Fijian: active transitive, passive transitive and transitive only. In
 6 order to express passive meaning, Fijian relies on the so-called passive transitive verbs
 7 used without the transitive inflectional marker (Geraghty 2008: 28-29). According to
 8 Biewer (2009: 366), Fijian employs several inflectional strategies to mark the passive
 9 on the verb, and according to Schütz (2014: 104) one of them is to add a stative marker
 10 to an active verb, which then has the goal as its subject, but without being able to add
 11 the agent in this construction:

12 ... there is not simple answer to the question: Does Fijian have a passive
 13 construction? If “passive” means simply “goal-focussed”, the answer is yes. If one
 14 insists that a passive sentence contain reference to the actor, the answer is no.

15 To sum up, all substrate languages have means of expressing the passive
 16 notionally, but formal means of doing so are present in the substrate of only four out
 17 of the five contact varieties we study. Table 1 summarises the means that are used in
 18 the various substrate languages by country.

	periphrastic	inflectional	other
Singapore	-	-	+
Hong Kong	-	-	+
India	+	-	-
Philippines	-	-	-
Fiji	-	+	-

19 **Table 1: Synopsis of grammatical passives in substrate languages**

20 If substrate influence were to play a significant role, one possibility is that the
 21 combination of passive with other verbal categories is avoided in SingE and HKE,
 22 where the dominant substrate languages mark passive outside the VP and where
 23 complex VPs may occur less frequently than in other WEs. We further expect
 24 ‘givenness’ to play a more important role in PhilE than in the other ESL varieties if
 25 substrate plays a role: according to Maratsos (1988: 133)

26 [e]ach sentence in Tagalog has to give at least one argument that is
 27 pragmatically given. If the sentence has a verb, the verb must have an
 28 agreement marker that marks the grammatical case (nominative, accusative,
 29 dative, or benefactive, typically) of this given argument.

1
2 These are possible outcomes of more subtle, probabilistic substrate influence than the
3 one evidenced in structural transfer that gives rise to constructions like the *kena-*
4 passive in SingE.

5 **2.5 Contact-induced phenomena and the passive**

6 In addition to substrate influence, processes of language acquisition may give rise to
7 features in the contact varieties which, in turn, may affect the realization of the *be-*
8 passive in the ESL varieties investigated here. A wide-spread phenomenon is final
9 consonant cluster reduction, for instance, which in turn may lead to the use of
10 unmarked past participles, as in the following example (retrieved from the *Global*
11 *Web-based English corpus*):³

12 (1) ... the land *was consider* worthless.
13 (<http://www.phuket101.net/2011/06/bang-tao-beach.html>)

14 This is said to be especially frequent in SingE and HKE, as expected from the isolating
15 nature of the substrate languages and from the fact that their substrate languages
16 tend to simplify consonant clusters (Deterding 2007: 17-19; Deterding *et al.* 2008).
17 Unmarked participles are likely to pose a problem for the automatic retrieval of
18 passives from syntactically annotated data (see section three).

19 Auxiliary deletion is another common feature found in various contact
20 varieties of English. Biewer (2015: 186) provides the following passive with auxiliary
21 deletion from her corpus of FijE (interview data):

22 (2) and in school when we speak Fijian we are we \emptyset always punished/
23 (Fiji/WI.txt)

24 This kind of feature is much less likely to occur in published academic writing than in
25 spoken interaction, however.

26 Previous research has also shown that some L2 learners of English tend to over-
27 passivise, that is, they extend the use of the passive to verbs that do not allow for a
28 passive in ENL varieties; these include unaccusative verbs that may or may not have
29 transitive counterparts in some contexts (*Tom ate / Tom arrived*) and unergatives,

³ GloWbE can be accessed at the following site <http://corpus.byu.edu/glowbe/>; for background information and a critical discussion of the advantages and shortcomings of this resource, see Davies and Fuchs (2015).

1 such as *Tom laughed* (Kondo 2005: 129). An example of overpassivization from learner
2 English is the following (from Kondo 2005: 156):

3 (3) The coin *was vanished* instantly.

4 On the whole, the structurally markedly different passive constructions that are due
5 to general processes of second language acquisition are quite likely to have been
6 edited out in published academic writing. Hundt et al. (2016) test for recall of
7 automatically retrieved transitive active and passive constructions and are able to
8 show that such characteristically ESL constructions are very infrequent in the
9 academic part of the ICE corpora.

10 **2.6 Research questions**

11 On the basis of previous research and the choice of varieties we investigate, we can
12 formulate the following research questions:

- 13 • Will we see a divide into first- and second-language varieties when it comes
14 to the internal factors predicting the choice of a passive over an active?
- 15 • More specifically, will we find regional differences in the role that animacy
16 of the subject plays as a predictor?
- 17 • For the second-language varieties, can we observe possible influence of
18 substrate languages or the process of second-language acquisition in the
19 effect size that the factors have in predicting voice alternation?

20 **3. Data and methodology**

21 **3.1 Parsed ICE corpora**

22 Our data come from the published academic writing section of the ICE corpora (see
23 Table 2). These comprise ten 2,000-word samples each from the humanities, social
24 sciences, natural sciences and technology, and thus a total of approximately 80,000
25 words per variety and 560,000 words for all Englishes investigated in this paper. For a
26 study of a frequent alternation as that between active and passive, such a relatively
27 small corpus yields ample evidence.

28

1

ICE-GB	Great Britain
ICE-US	United States of America
ICE-SIN	Singapore
ICE-HK	Hong Kong
ICE-IND	India
ICE-PHI	The Philippines
ICE-FJ	Fiji

2

Table 2: List of ICE components (abbreviations) included in the study

3 The ICE corpora were syntactically annotated with the probabilistic dependency
4 parser (Pro3Gres; Schneider 2008). The annotation process includes part-of-speech
5 tagging, chunking and parsing and makes automatic retrieval of passives and active
6 transitives possible (see next section).

7 3.2 Data retrieval

8 Hundt et al. (2016) comment in detail on the automatic retrieval of central *be*-passives
9 and active transitives. We make use of the same approach to retrieve data from the
10 seven ICE components. We then coded sets of 200 randomly sampled instances per
11 variety, 100 passives and actives, each, and equal numbers of actives and passives per
12 subdiscipline. We manually excluded false positives such as the following from our
13 data.

14 (4) ... more Fijian businesses would *bite* dust if Government did not act now
15 (ICE-FJ W2A 011)

16 (5) Therefore it can be assumed that those who *enter* a university have chosen
17 to do so ... (ICE-SL W2A 004)

18 (6) One of the problems I *wish to address* is the degree to which Frankish uncial
19 in the late eighth and the ninth centuries is indeed artificial rather than
20 natural. (ICE-GB W2A 008)

21 Likewise, only those central *be*-passives were included in our investigation that had
22 active transitive counterparts, excluding for instance adjectival passives as in (7):

23 (7) Mesozoic isolated platforms *are well represented* in the Tethyan region.
24 (ICE-GB W2A 023)

25 Having excluded non-relevant material from our randomly sampled concordances, we
26 initially ended up with a total of $N=1285$ variants (roughly equal amounts of actives
27 and passives) that we coded for various internal factors.

1 3.3 Factors coded for

2 *Complexity of the verb phrase*

3 We coded the complexity of the VP in terms of TAM combinations with voice,
4 distinguishing between simple present and past, progressive present and past,
5 present and past perfect as well as combinations of central modals with active or
6 passive. As it turned out, the factor was only weakly significant as long as individual
7 patterns were coded separately (e.g. present and past progressive). We therefore
8 decided to group complex VPs (perfect, progressive and modal) together and contrast
9 these with simple (both present and past) VPs.

10 *Weight*

11 The concept of weight has been defined as number of syllables, words, nodes and
12 phrasal nodes (Rosenbach 2007; Hawkins 2004). As Wasow & Arnold (2003) point out,
13 however, it is still not clear whether the different ways of measuring constituency
14 weight yield different results (see Seoane 2009: 370-371). For the purposes of our
15 study, we measure constituent length in number of words; for constituents that are
16 longer than four words, we use brackets (4-9 coded as '4', 5-10 as '5', etc.). This means
17 that in some examples, the subject may actually be slightly longer in total number of
18 words than the object, even though in our analysis they show up as having the same
19 length in terms of the bracket they occur in. An example is given in (8), where the
20 subject amounts to 9 words whereas the object comprises only 5 words, with both
21 falling into our coding bracket '4', i.e. constituents between 4 and 9 words in length:

22 (8) ... [the high amplitude specular signal from the 0 probe]_S emphasises [the
23 surface of the defect]_O. (ICE-GB W2A 031)

24 Also, we counted the number of words of the syntactic unit functioning as the subject
25 and object/agent rather than the semantic subject. In the following example (9), it is
26 the relative pronoun that is the syntactic subject (1 word) of the active transitive verb
27 *change*, even though its antecedent, and thus the semantic subject (*massive*
28 *upheavals ...*), is much longer.

29 (9) World War II produced massive upheavals beyond the preemptive political
30 and military events that *changed* the future course of history. (ICE-US W2A
31 013)

32

1 *Animacy/Semantics*

2 The factor ‘animacy’ is not a binary one but a gradient (see section 2.2).⁴ Nonetheless,
3 we coded for animate subjects/objects and inanimate ones, adding a third category
4 (unclear) only if animacy could not be determined on the basis of the context.⁵ This
5 means that the factor ‘animate’ in our analysis includes a fairly broad range of
6 semantic concepts, such as prototypical animate and human NPs, e.g. personal
7 pronouns as in (10) and the authors’ names in (11), but also collective nouns such as
8 *government* and personifications such as *Rome* in examples (12) and (13),
9 respectively.

10 (10) what I found was that almost all his poetry is filled with expressions of his
11 need for and love of other people. (ICE-GB W2A 003)

12 (11) *Wilson & Henderson* (1966) described *U. ambiguus* and *P. allii* as the two
13 species in the U.K. with *P. porri* and *P. mixta* (on chives) as synonyms of *P.*
14 *allii* (ICE-GB W2A 028)

15 (12) The Griffiths Report was not received with great enthusiasm by *the*
16 *Conservative government* ... (ICE-GB W2A 013)

17 (13) ... and in the stress of war *Rome* conceded what they had sought. (ICE-GB
18 W2A 001)

19 Examples of constituents whose semantics are unclear are given in (14) and (15):

20 (14) These steps are iterated until an acceptable solution is reached. (ICE-GB
21 W2a 016)

22 (15) So that neither *the user or system* is overwhelmed by large result sets, the
23 size of result sets is limited to 100 items ... (ICE-US W2A 038)

24

25 The agent could be an abstract process in example (14) or a human agent, whereas
26 the combination of an animate with a clearly inanimate NP in (15) renders the subject
27 constituent neutral, so to speak, with respect to the factor ‘animacy’. Cases of unclear
28 animacy coding are discarded in the subsequent analyses.

29 In short passives, where the semantics of the agent could be inferred from the
30 context, animacy was coded for as outlined above. In example (16), for instance, the
31 inferred object/agent is coded as ‘animate’, while in example (17), the inferred agent
32 of the passive phrase is coded as ‘inanimate’.

⁴ Zaenen et al. (2004), for instance, distinguish eleven categories, including two for instances where the annotator was unsure of the coding. This indicates that there is no universal solution to the problem of coding for the factor ‘animacy’.

⁵ These had to be excluded in the mixed-effects modelling because of model convergence issues, leaving us with 1,168 items.

1 (16) For example, thermotropic responses of seedlings *have* occasionally *been*
2 *noted* (Aletsee 1962a), [...]. (ICE-GB W2A 025)

3 (17) Leaching is a continuous process in the dunes because the sands do not
4 overlie a bedrock *which is being gradually weathered*, as do most terrestrial
5 soils, but if anything are accreting more depth of sand. (ICE-GB W2A 022)

6 *Definiteness/Givenness*

7 Finally, the factor ‘givenness’ is not really a binary one, either (see section 2.2). There
8 are various ways of coding for ‘givenness’ (for a discussion, see Dreschler 2015: 83-
9 86). We decided to approximate this factor by coding both active objects and passive
10 subjects as either ‘pronominal’, ‘definite’, or ‘other’. Due to data sparseness,
11 pronominal and definite constituents were then conflated under the label ‘definite’,
12 all other instances were coded as ‘indefinite’. Among definite noun phrases, we also
13 included instances with proper names:

14 (18) *Mrs. Roca* ... expressed her frustration over the intricate process she had to
15 go through to face her husband about the separation. (ICE-PHIL, W2A 015)

16

17 **3.4 Statistical modelling**

18 In order to tease apart influence of the predictors on voice alternation across the
19 seven varieties of English, we make use of two multivariate techniques. First, we
20 model the data using a tree and forest approach as first advocated by Tagliamonte &
21 Baayen (2012). Both are variants of permutation testing that do not assume a certain
22 distribution of the data but build a model by resampling from the input. The random
23 forest analysis provides information on overall variable importance and the single
24 conditional inference tree (ctree) allows us to tap into and visualise possible
25 interactions between predictor variables.⁶ The ctree splits the data recursively into
26 smaller subsets according to those predictors that co-vary most strongly with the
27 outcome. For each binary split, the data is inspected for the predictor that best
28 preserves the homogeneity of each split (e.g. all actives versus all passives) at the
29 customary significance level of $\alpha = 0.05$. The splitting process is then repeated until no
30 further splits can significantly increase the subsets’ homogeneity with regard to the

⁶ We used the party package (Strobl, Hothorn and Zeileis 2009) and the partykit (Hothorn et al. 2006) as implemented in R, respectively, to fit the conditional random forest and the ctree.

1 outcome variable. The bottom-most barplots provide information on the observed
2 proportions of outcome variants in that particular split. Single trees have the
3 disadvantage that they very much hinge on the dataset at hand and are hence subject
4 to a high degree of variability. Conditional random forests resample over a predefined
5 number of trees using a conditional permutation scheme and are thus particularly
6 robust to, for instance, predictor collinearity. In addition to predictor rankings
7 obtained through the random forest, and in order to gain a more robust insight into
8 the predictors' variability across our set of geographically dispersed varieties of
9 English, our second approach uses mixed-effects modelling.⁷ This allows us to
10 evaluate the significance of individual predictor variables in the voice alternation
11 (Hosmer and Lemeshow 2000; Pinheiro and Bates 2000). A mixed-effects model
12 makes adjustments to the model's predictions from the fixed effects by including
13 random effects in the modelling process. Random effects account for idiosyncratic
14 variation by group that is specific to the dataset, such as lexical items or texts sampled.
15 By using mixed-effects modelling, we are able to generalise beyond the particular
16 dataset at hand to, for instance, all texts of academic English. In particular, it allows us
17 to look into possible interaction with the predictor 'variety', which our sampling
18 precluded from emerging as an important effect in the random forest analysis.

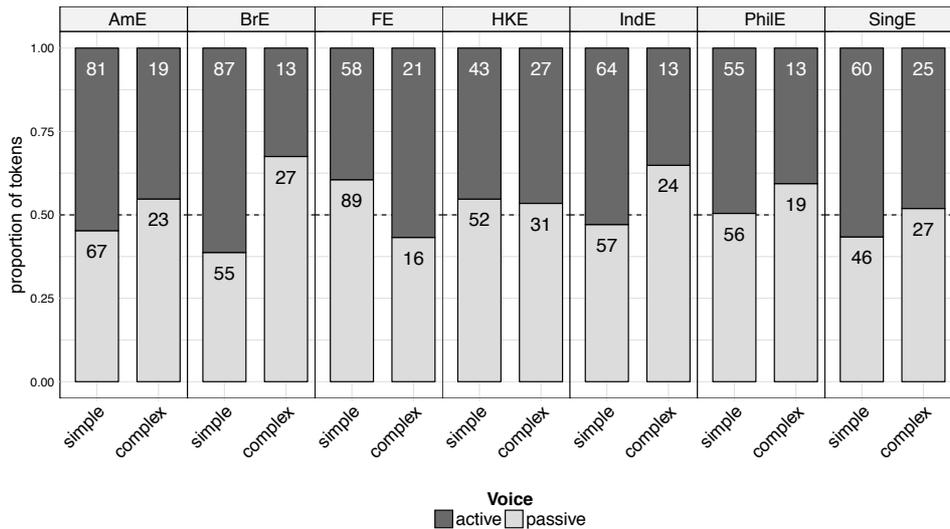
19 **4. Results**

20 **4.1 Descriptive statistics**

21 *Complexity of the verb phrase*

22 The proportional distributions in Figure 2 do not confirm our hypothesis of substrate
23 influence, which would have predicted that passives are avoided in complex VPs in
24 varieties where the main substrate marks voice outside the VP (i.e. SingE and HKE).
25 Instead, Figure 2 shows that passives are proportionally preferred over actives in
26 complex VPs in academic writing, with the exception of FijE. In simple VPs, actives tend
27 to be preferred over passives, with the exception of FijE and HKE, which prefer the
28 passive, and PhilE, which shows an even distribution.

⁷ We make use of the `glmer()` function from the `lme4` package in R (Bates et al. 2015) for this.

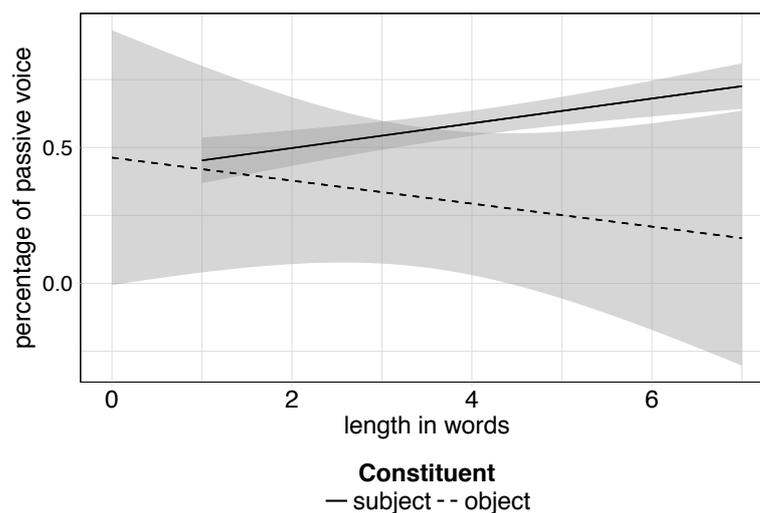


1

2 **Figure 2. Proportional distribution of simple and complex VPs by voice variant. Raw numbers provided in each**
 3 **bar.**

4 *Weight*

5 Figure 3 shows the smoothed conditional means of the proportion of passives (y-axis)
 6 by increasing length (x-axis) of the subject (solid line) and the object (dashed line): the
 7 proportion of passives increases with increase in subject-length and the proportion of
 8 actives increases with increase in object-length, in accordance with the principle of
 9 end-weight.



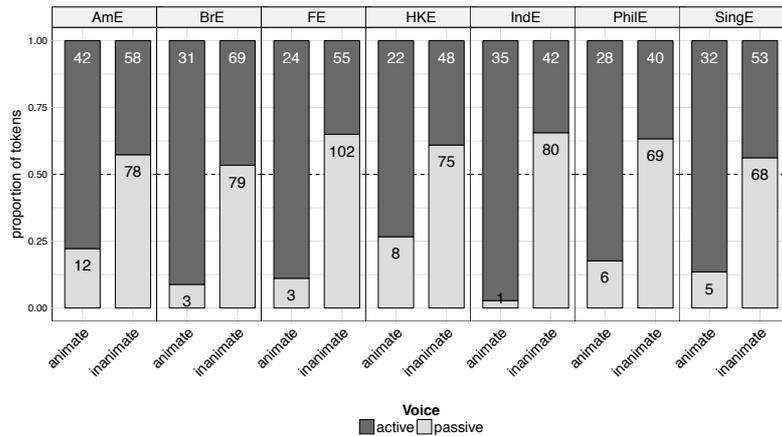
10

11 **Figure 3. Smoothed conditional means of the proportion of passive variants by increasing object length (dashed**
 12 **line) and subject length (solid line). Note that objects can have 0 length when they are not present in**
 13 **the construction.**

14 *Animacy/Semantics*

15 Regarding the proportional distribution of SemanticsOfSubject, actives prefer animate

1 subjects whereas inanimate subjects dominate in the passive (see Figure 4).

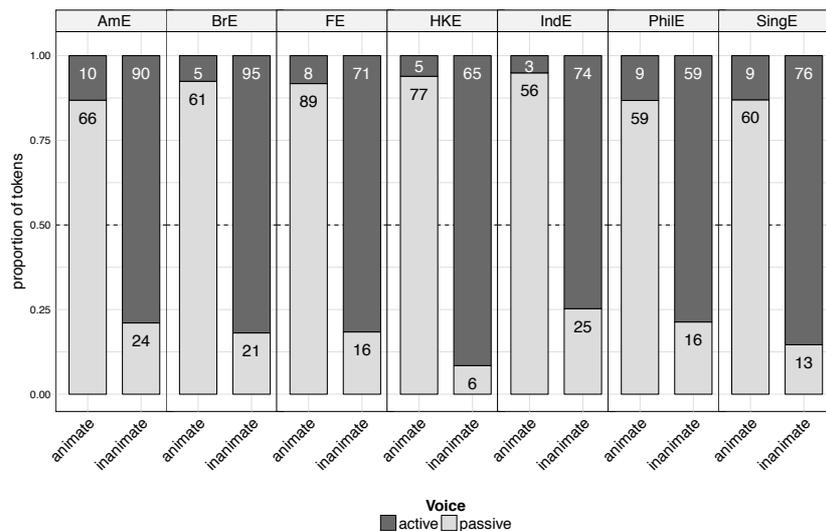


2

3 **Figure 4. Proportional distribution of animate and inanimate subjects by active and passive voice variants.**

4 With respect to SemanticsOfObject, the situation is reversed with actives preferring
 5 inanimate objects and passives showing an overall preference for animate objects (see
 6 Figure 5).

7

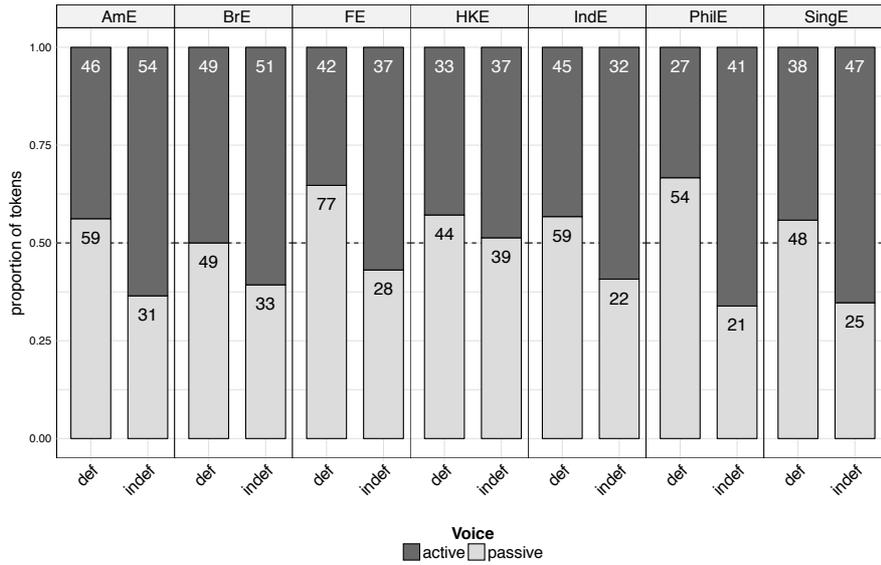


8
9

Figure 5. Proportional distribution of animate and inanimate objects by active and passive voice variants.

10 *Definiteness/Givenness*

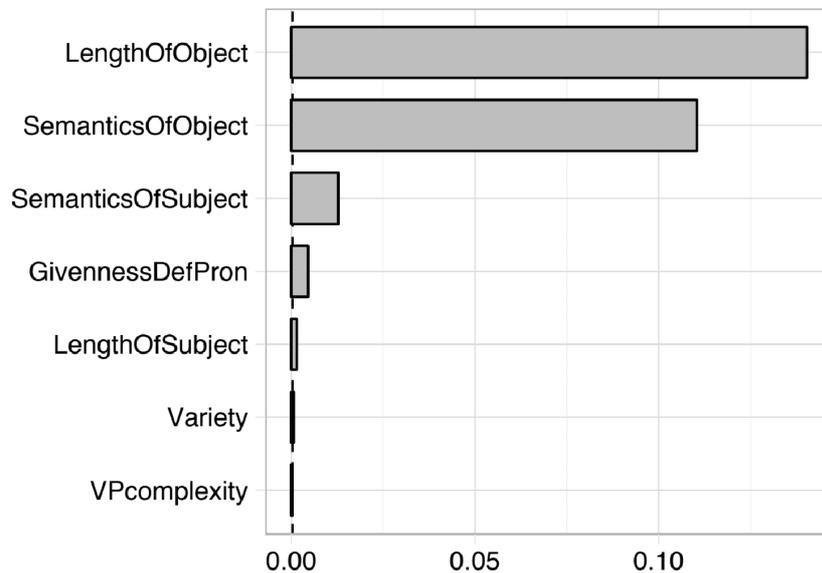
11 The general trend across varieties is for definite patients to be more likely in
 12 passives, where they function as subjects, than in the active voice. Indefinite
 13 patients, on the other hand, are more likely in actives, where they function as
 14 objects. Definiteness/Givenness has thus the expected effect reported in the
 15 literature.



1
2 **Figure 6. Proportional distribution of definite and indefinite patients by active and passive voice variants.**

3 **4.2 Forest and tree analysis**

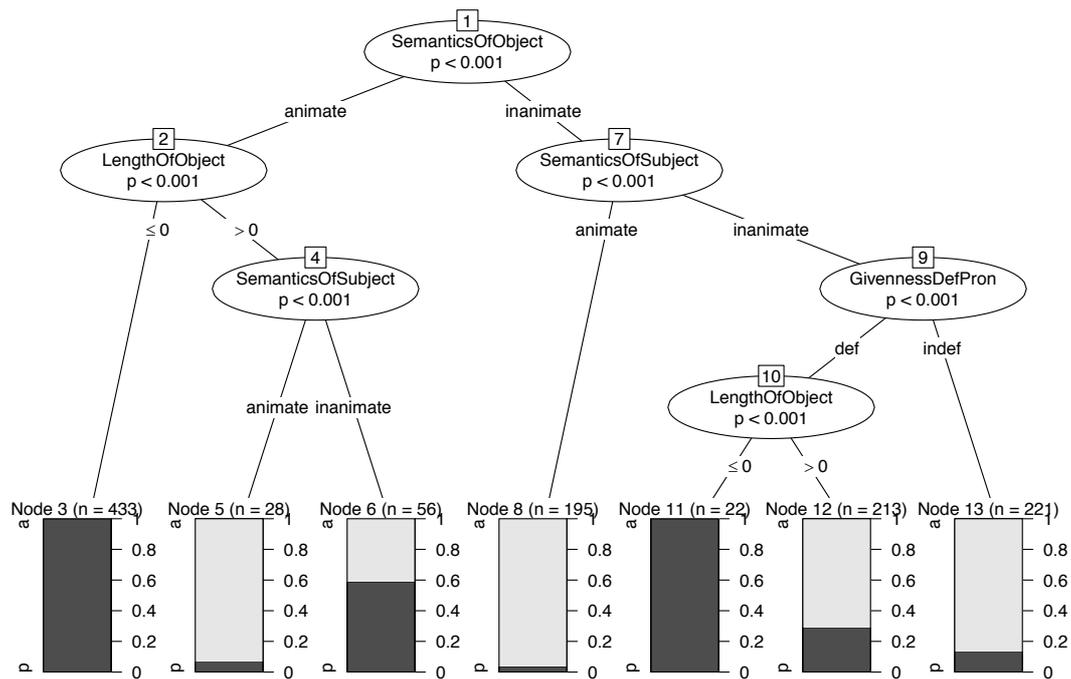
4 Figure 7 ranks the importance of the predictor variables. It shows that length of object
 5 (related to the principle of end-weight) is the most important predictor for the choice
 6 of a passive over an active VP, followed by animacy (notably the semantics of the
 7 object). All other factors are less important.



8
9 **Figure 7. Overall variable importance (random forest analysis) for voice alternation in academic English⁸**

⁸ We set mtry=3 and ntree=2000, following recommendations by Strobl (p.c.). Somers2 Dxy returns a prediction accuracy of 0.951 and a C-index of 0.975, which is above the level of 0.8 recommended e.g. in Tagliamonte & Baayen (2012: 156).

1 The output of a conditional inference tree is shown in Figure 8.⁹



2

3 **Figure 8. Conditional inference tree of passive (p) over active (a) choice in academic English.**

4 The ctree returns SEMANTICSOFOBJECT as the most predictive factor (Node 1). Inanimate
 5 objects are more likely to be realised in the active variant, while animate objects (and
 6 unclear cases) favour the passive. Proceeding down the left side of the tree, we see
 7 that among animate objects and unclear cases, the passive variant is categorically
 8 chosen if the object is inexistent (Node 2). This is not surprising given that animate
 9 agents are hardly ever expressed in passives and often inferred from the context (see
 10 example (16)). In cases where the object/agent is present, inanimate subjects have a
 11 stronger preference for the passive than animate subjects (Node 4).

12 Moving down the right side of the tree where the object is inanimate (from
 13 Node 1), we see again a split by semantics of the subject (Node 7): In the case of
 14 inanimate objects, the active variant is preferred if the subject is animate (in other
 15 words, both subject and object are expressed in active transitive clauses) (Node 8).

⁹ The stability of the classification was confirmed through a second tree fitted with a different random seed. Classification accuracy of the conditional inference tree is 89.4%, i.e. significantly above the baseline of 50.4%; the C-statistic (0.894) is also above the recommended level of 0.8. Note that only splits with a *p*-value of 0.05 or lower were allowed when building the tree.

1 When both subject and object are inanimate, the definiteness of the syntactic object
2 in the active and the syntactic subject in the passive clause as well as length of the
3 object become relevant (Nodes 9 and 10). Passive is the categorical choice if the
4 syntactic subject is definite or pronominal and the object/*by*-agent is not explicitly
5 expressed (passive is the only viable option in those cases) (Node 11). If the object is
6 expressed, however, actives are the preferred option (Node 12). If both object and
7 subject are inanimate the active is preferred with indefinite subjects (Node 13).
8 Neither 'variety' nor 'subdiscipline' show up as significant factors most probably due
9 to the sampling process.

10 **4.3 Mixed-effects model**

11 Our initial model included all factors (Section 3.3) as fixed effects in interaction with
12 VARIETY (no other higher order interactions were considered). The numeric factors
13 LENGTHOFOBJECT and LENGTHOFSUBJECT were standardised by two standard deviations
14 and centered around the mean in order to reduce potential covariation among the
15 numeric and other predictors and to create a predictor with an effect size on a scale
16 comparable to that of a binary predictor (see Gelman 2008). VARIETY was coded using
17 sum contrasts whereby the proportion of responses for each level is compared against
18 the grand mean across all levels (see Menard 2010: 97). Random effects included a
19 simple random intercept for File in order to account for idiosyncracies in the texts
20 sampled to represent academic writing.

21 Model selection followed the backwards elimination procedure outlined by
22 Zuur et al. (2009: 120–122). Because of convergence issues with the full set of
23 interactions, we additionally fitted a fixed-effect model in order to identify
24 interactions that did not significantly improve model fit and cross-tabulated the data
25 to find sparse cells. The exclusion of the most negligible predictor (SEMANTICSOFSUBJECT)
26 from any interactions made model convergence possible. The predicted outcome of
27 the model was the log odds of the passive variant. Next, we identified and excluded
28 any other interaction terms that did not significantly improve model fit.

29 The final model (19) includes a by-file random intercept, as well as an
30 interaction of VARIETY and SEMANTICSOFOBJECT, GIVENNESS, and LENGTHOFSUBJECT. None of
31 the other initial interactions with VARIETY turned out to be significant.

1 (19) Passive model 1; Response = {active, passive}
 2 Response ~ (1 | FILEID) + SEMANTICSOFSUBJECT + VP COMPLEXITY + LENGTHOF OBJECT
 3 + VARIETY * (SEMANTICSOF OBJECT + GIVENNESS + LENGTHOF SUBJECT)

4 Classification accuracy of the model is 88.5% which is significantly better than the
 5 baseline of 50.43% when always choosing the most frequent (passive) variant (p_{binom}
 6 < 0.001). Summary statistics further give a very good index of concordance $C = 0.946$,
 7 which indicates a good model fit (see Baayen 2008). The condition index $\kappa=10.9$
 8 indicates existent but not harmful collinearity (Baayen 2008: 182). The variance
 9 inflation factor (VIF) for each of the factors points out that much of the estimated
 10 variance of all higher order interactions with VARIETY is associated with the
 11 corresponding main effect and with other interaction effects with VARIETY. We will
 12 therefore exercise extra caution when interpreting these results.

13 To evaluate the model, we randomly divided our dataset 100 times into a
 14 training and a test set. Next, we fitted the model to each training set and calculated
 15 the predictions for the corresponding test set. The accuracy and index of concordance
 16 C of each model was then measured and compared to the original model. Mean
 17 accuracy of the 100 models was 86.8% which corresponds to a drop of 1.7% from the
 18 original model. C -statistic was an excellent 0.946.

19 *Main effects*

20 The coefficients of the main factors in the model are summarised in Table 3.

Factor	β	SE	p
(INTERCEPT)	-4.25548	0.45754	<0.001 ***
SEMANTICSOF SUBJECT			
animate \Rightarrow inanimate	2.64365	0.31323	<0.001 ***
SEMANTICSOF OBJECT			
inanimate \Rightarrow animate	3.31353	0.38393	<0.001 ***
GIVENNESS/DEFINITENESS			
indef \Rightarrow def	1.10909	0.22423	<0.001 ***
LENGTHOF OBJECT	-2.22938	0.28321	<0.001 ***
LENGTHOF SUBJECT	0.91655	0.29722	<0.00204 **

21 Table 3: Main effects of individual factors in the model. Predictions are for passive voice. Only significant
 22 factors shown.

23 The column labelled β indicates the estimates of the coefficients on a logit-scale.
 24 Positive values signal a preference for passive (the predicted outcome), negative
 25 values a preference for the active voice. SE specifies standard errors. The results of

1 the logistic regression analysis can be summarised as follows:

2 First, the main effects in the model have largely the predicted influence on the
3 choice of voice variant given previous literature. The probability of the passive
4 increases if the subject is inanimate instead of animate, if the object is animate instead
5 of inanimate, if the VP is simple instead of complex (not a significant factor) and with
6 every unit increase in the length of the subject. In other words, passives tend to be
7 used with an inanimate long subject and animate object/agent. The likelihood of a
8 passive also increases if the subject of the passive voice or the object of the active
9 voice is definite, that is, definite constituents (e.g. pronouns) are more likely to be
10 used in subject position in passives than in object position in active voice. This follows
11 naturally from the ordering preference in the human processing system in that
12 speakers tend to use definite/given constituents before new information (see Wasow
13 2002: 30 and literature cited therein; also Ransom 1979). The probability of a passive
14 further decreases with every unit increase in the length of the syntactic object, i.e. the
15 agent in the passive voice. Table 3 only reports significant factors; in other words, VP
16 complexity turns out not to be a significant predictor. Cross-varietal differences do not
17 emerge as a main effect in the overall choice of passive vs. active because equal
18 numbers of passives and actives was sampled per variety and subdiscipline. The next
19 section therefore specifically looks at possible interaction effects with VARIETY.

20 *Interaction terms*

21 Table 4 reports all significant interaction terms between VARIETY and the language-
22 internal factors SEMANTICSOFOBJECT, GIVENNESS and LENGTHOFSUBJECT. Note that none of
23 the other interaction terms contribute significantly to the model fit. If the coefficient
24 estimates of a main predictor (for instance, SEMANTICSOFOBJECT) and its interaction
25 term (VARIETY : SEMANTICSOFOBJECT) have the same +/- sign, the effect is stronger in that
26 specific variety (compared to all other varieties). If they have opposite signs, the effect
27 of that factor is weaker in that specific variety (or possibly even reversed). Predictions
28 are for the passive variant.

29

1

Factor	β	SE	p
VARIETY : SEMANTICSOFOBJECT			
HKE + animate	3.37523	1.49674	0.02414 *
VARIETY : GIVENNESS/DEFINITENESS			
HKE + def	-1.48428	0.68535	0.03033 *
PhilE + def	1.12773	0.54831	0.03971 *
VARIETY : LENGTHOFSUBJECT			
HKE	3.38519	1.40018	0.01562 *

2

Table 4: Significant interaction effects between Variety and language-internal factors in the model. Predictions are for passive voice.

3

4 Zooming in, the interaction effects in the model indicate that academics from Hong
5 Kong deviate from the global average with regard to the effect of all three language-
6 internal predictors, and those from the Philippines with regard to only one.

7

- Animate objects are more likely to be used in passives in HKE academic writing;
8 the effect of SEMANTICSOFOBJECT is thus stronger in that variety (see Figure 3).

9

- At the same time, the effect of GIVENNESS/DEFINITENESS is weaker in HKE, that is,
10 if the constituent is definite, the likelihood of passive voice is not as strong as
11 in the other varieties (see Figure 4).

12

- The effect of GIVENNESS/DEFINITENESS is stronger in the academic writing of
13 researchers from the Philippines, i.e. they are more likely to use passive voice
14 if the constituent is definite (see Figure 4).

15

- The effect of LENGTHOFSUBJECT is stronger in academic writing in HKE compared
16 to the global average: academics from Hong Kong are more likely than
17 academics elsewhere to use a passive when the syntactic subject increases in
18 length (see Figure 5).

1 Figures 9 to 11 report the probability of passive voice given the predictor's level, for
 2 instance animate vs. inanimate (y-axis), across the seven varieties (x-axis). Fitted
 3 effects were plotted using the effects package in R (Fox 2003).

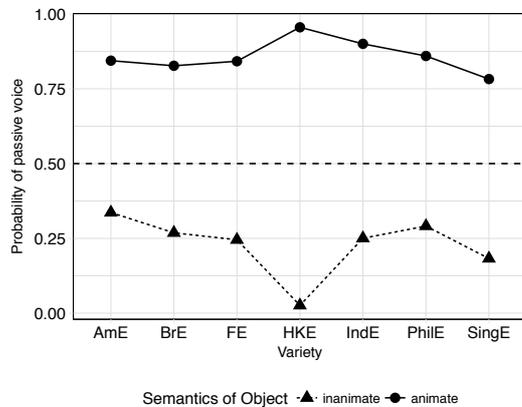


Figure 9: Passive voice and Semantics of object

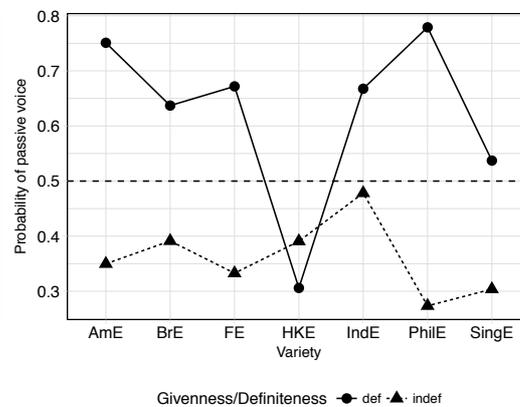


Figure 10: Passive voice and Givenness

4 Figure 9 shows that the passive is more likely with animate objects and active voice
 5 with inanimate ones. Moreover, the difference in the effect is strongest in HKE. The
 6 graphs in Figure 10 reveal that the probability of passive is generally higher with
 7 definite than indefinite constituents: Givenness/Definiteness has a reverse effect in
 8 HKE and a stronger effect in PhilE academic writing. Finally, with respect to the effect
 9 that the length of the subject has on the choice of a passive over an active VP, only
 10 academic texts from Hong Kong show a significant effect of this predictor, but none
 11 that could be easily explained by substrate influence (Figure 11).

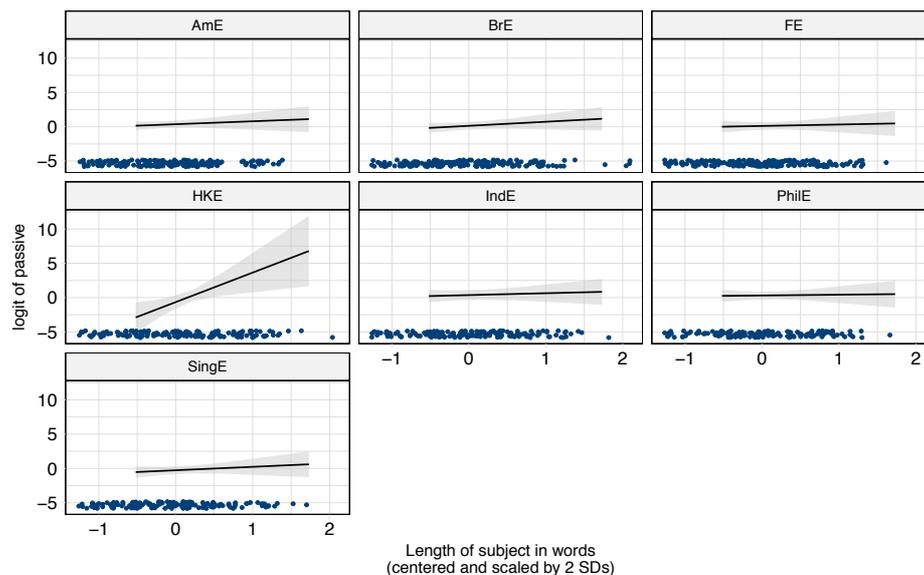


Figure 11: Effect of increase in subject length by variety.

12
 13

1 **5. Summary and discussion**

2 The analyses show that the choice of a passive over an active VP in academic writing
3 is only partially subject to regional variation. Importantly, there is no straightforward
4 difference between ENL and ESL varieties (i.e. types of Englishes) when it comes to the
5 factors predicting choice of a passive. It is only in academic writing from Hong Kong
6 and the Philippines that we can observe very subtle cross-lectal variability which is to
7 be found in the effect size of LENGTHOFSUBJECT, SEMANTICSOBJECT and GIVENNESS.
8 Overall then, the use of passive voice in academic writing is fairly homogenous across
9 different varieties of English, be they norm-providing ENL varieties or ESL varieties.
10 How do these results connect with our initial research questions?

11 Previous research into other morphosyntactic alternations found that
12 ‘animacy’ was often a factor that had variable influence across national varieties (see
13 section 2.2). ‘Animacy’ is relevant for voice alternation in two (related) respects.
14 Firstly, the passive is a construction that allows the subject to be inanimate, contrary
15 to the usual animacy hierarchy.¹⁰ Secondly, the construction also allows the demotion
16 of the agent (typically animate) to the object relation (either overtly expressed as a
17 *by-agent* or to be inferred from the context). While this is a global constraint that holds
18 in academic writing across the varieties investigated, our data show that animacy is a
19 particularly strong factor in Hong Kong academic writing, which prefers passives with
20 animate (underlying) agents. Representative examples that imply an animate, human
21 agent which, however, typically remains unexpressed, are the following:

- 22 (20) *It is believed* that the internal cohesion within the relatively self-contained
23 or isolated construction activities will be impaired if adversarial proceedings
24 are employed to settle disputes. (ICE-HK, W2A-014)
25 (21) Cases of amoebic meningoencephalitis due to *Naegleria fowleri* *have also*
26 *been seen* locally. (ICE-HK, W2A-021)
27 (22) What *was forgotten* at the moment or what *was repressed* in the past
28 returns to consciousness in dreams, helping men to re-gain the balance of
29 a healthy mind. (ICE-HK, W2A-002).

¹⁰ Note that occasionally, inanimate subjects are used in active transitives even though the patient in object position is animate, as in the following example: “*Section 29* did empower *local authorities* to promote the welfare of persons who are blind, deaf or dumb, ... ” (ICE-GB W2A 013).

1 Givenness turned out to have a weaker effect in HKE than the other varieties. The
2 following are typical HKE examples of passives with non-specific subjects:

3 (23) *Certain types of building components* are required to be tested for their
4 safety and serviceability. (ICE-HK, W2A-036)

5 (24) *Occasional cases of fungal abscesses* have also been seen. (ICE-HK, W2A-
6 021)

7 (25) *Other arts of Chinese* were also combined to the lantern design, ... (ICE-HK,
8 W2A-008)

9 The reverse holds for PhilE, where passivised subjects frequently contain a
10 demonstrative pronoun, a clear marker of 'givenness', or are otherwise clearly 'given'
11 in the preceding context (over and beyond being 'definite'):

12 (26) *These cut portions* are left for sometime at ambient condition to cure
13 wounds naturally. (ICE-PHI, W2A-033)

14 (27) At the very outset, *this leitmotif* has to be grasped well because many
15 concepts of Asian philosophy cannot be understood without it. (ICE-PHI,
16 W2A-009)

17 (28) It is true that an object can require another to perform some task, but, *the*
18 *actual manipulation of the required data* or service is done in the called
19 object. (ICE-PHI, W2A-038)

20 That givenness turns out to play such an important role in PhilE academic writing fits
21 in well with the dominant substrate, Tagalog (see section 2.4).

22 With respect to the factor 'weight', the default in English (and other languages)
23 is to have short subjects and long objects. A good example of a prototypical
24 combination of short subject and heavy object is given in (29), where the reverse
25 achieved by a passive would be highly unlikely.

26 (29) I began the process of questioning and reexamining the possibility of
27 comparatively studying a select number of successful nonprofit
28 organizations and their leadership to determine what strengths would be
29 revealed regarding their success. (ICE-US, W2A-011)

30 If both subject and object/agent are animate, length is the deciding factor:

31 (30) Now in her twilight years, she was a nervous recluse, living off the charity
32 of family and friends, and *she was only visited regularly by a hired reader*.
33 (ICE-US, W2A-007)

34
35 In passives, long subjects typically occur without agent *by*-phrases, for instance in
36 contexts where the agent is unknown.

1 (31) *The two earliest, the Aberdeen fragment written in rustic capitals and the St*
2 *Gall Vulgate in half-uncials*, may have been written at a time when the
3 hierarchy of appropriate scripts was being worked out. (ICE-GB, W2A-008)

4
5 Note that the length of the subject in (31) is actually due to appositions that follow
6 the subject NP. Clausal expansion of the subject NP also contributes the bulk of
7 material in the two extremely ‘weighty’ passive subjects in the HKE data (examples
8 (32) and (33)), but long, internally complex phrasal NPs are also attested (34):

9 (32) *Traditionally, site control tests and sampling, such as slump and compacting*
10 *factor tests to CS1; flow of fresh concrete to BS 1881: Part 105; and making*
11 *of concrete cubes to CS1, etc*, were performed by the contractor’s site staff.
12 (ICE-HK, W2A-036)

13 (33) *Well-designed population-based studies that aim to access the long-term*
14 *significance of these new recommendation to diagnosis categories of*
15 *glucose intolerance in Asian populations* are now needed. (ICE-HK, W2A-
16 028)

17 (34) *The angiotensinogen M235T (TT) genotype and its possible interaction with*
18 *the angiotensin converting enzyme deletion/insertion polymorphisms*
19 *(DI/DD)* have also been reported in Chinese diabetic patients who have
20 increased albuminuria. (ICE-HK, W2A-024)

21 That long passive subjects should occur with a greater-than-average frequency (see
22 Figure 11) does not find an easy language-internal explanation, however. On the
23 contrary, if substrate were to play a role, we would expect to find exactly the opposite,
24 i.e. relatively short sentences and phrases and no post-modifying clauses within NPs.
25 Our data come exclusively from academic texts, though, and the explanation for the
26 high incidence of long passive subjects in academic writing from Hong Kong might well
27 have to be attributed to the conscious attempt at emulating a western academic
28 writing style resulting in an over-use of constructions that are avoided by researchers
29 with English as their first language (see also Gunn 2017 on stylistic effects in Chinese
30 writing as a result of translation practices).

31 Finally, our study did not produce a regional difference in the effect size for
32 the factor ‘complexity of the verb phrase’. This means that, contrary to our hypothesis,
33 there is no statistically significant trend in ESL academic writing to avoid complex
34 passive VPs, any more than would also be the case in writing produced by academics
35 with English as their first language.

1 **Acknowledgements**

2 TEXT

3 **References**

4 *Corpora*

ICE *International Corpus of English*

5

6 *Bibliography*

7 AUTHORS. forthcoming.

8 Baayen, Harald. 2008. *Analyzing linguistic data: A practical introduction to statistics*
9 *using R*. Cambridge: Cambridge University Press.

10 Bao, Zhiming and Lionel Wee. 1999. The passive in Singapore English. *World Englishes*
11 18(1): 1-11.

12 Bates, Douglas, Martin Mächler, Benjamin Bolker, and Steve Walker. 2015. Fitting
13 Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software* 67(1): 1-
14 48.

15 Behagel, Otto. 1909. Beziehungen Zwischen Umfang Und Reihenfolge von
16 Satzgliedern. *Indogermanische Forschungen* 25: 110–142.

17 Behagel, Otto. 1930. Von deutscher Wortstellung. *Zeitschrift für Deutschkunde* 44: 81-
18 89.

19 Biber, Douglas and Edward Finegan. 1989. Drift and evolution of English style: a history
20 of three genres. *Language* 65: 487–517.

21 Biewer, Carolin. 2009. Passive constructions in Fiji English: A corpus-based study. In
22 Andreas H. Jucker, Daniel Schreier and Marianne Hundt, eds. *Corpora: Pragmatics*
23 *and Discourse*. Amsterdam: Rodopi, 361–377.

24 Biewer, Carolin. 2015. *South Pacific Englishes. A Sociolinguistic and Morphosyntactic*
25 *Profile of Fiji English, Samoan English and Cook Islands English*. Amsterdam:
26 Benjamins.

27 Bresnan, Joan and Jennifer Hay. 2008. Gradient grammar: An effect of animacy on the
28 syntax of *give* in New Zealand and American English. *Lingua* 118: 245–259.

29 Bresnan, Joan and Marilyn Ford. 2010. Predicting syntax: Processing dative
30 constructions in American and Australian varieties of English." *Language* 86(1):
31 168–213.

32 Davies, Mark and Robert Fuchs. 2015. Expanding horizons in the study of World
33 Englishes with the 1.9 billion word Global Web-based English Corpus (GloWbE).
34 *English World-Wide* 36(1): 1–28.

35 Deterding, David. 2007. *Singapore English*. Edinburgh: Edinburgh University Press.

36 Deterding, David, Jenny Wong and Andy Kirkpatrick. 2008. The pronunciation of Hong
37 Kong English. *English World-Wide* 29(2): 148–175.

38 Dreschler, Gea. 2015. *Passives and the loss of verb second: A study of syntactic and*
39 *information-structural factors* (LOT dissertation series 402). Utrecht: LOT.

40 Fox, John. 2003. Effect Displays in R for Generalised Linear Models. *Journal of*
41 *Statistical Software* 8(15): 1–27. URL <http://www.jstatsoft.org/v08/i15/>.

42 Gelman, Andrew. 2008. Scaling regression inputs by dividing by two standard
43 deviations. *Statistics in Medicine* 27(15): 2865–2873.

- 1 Gelman, Andrew and Jennifer Hill. 2007. *Data analysis using regression and*
2 *multilevel/hierarchical models*. Cambridge: Cambridge University Press.
- 3 Geraghty, Paul. 2008. *Fijian*. Victoria: Lonely Planet Publications.
- 4 Geraghty, Paul, France Mugler and Jan Tent, eds. 2006. *Macquarie Dictionary of*
5 *English for the Fiji Islands*. Macquarie University: The Macquarie Library.
- 6 Gunn, Edward. 2017. Westernization of Chinese Grammar. In Rint Sybesma, Wolfgang
7 Behr, Yuego Gu, Zev Handel and C.-T. James Huang, eds. *Encyclopedia of Chinese*
8 *Language and Linguistics*, Consulted online on 20 July 2017
9 <http://dx.doi.org/10.1163/2210-7363_ecll_COM_00000450>
- 10 Hawkins, John A. 2004. *Efficiency and Complexity in Grammars*. Oxford: Oxford
11 University Press.
- 12 Hinrichs, Lars and Benedikt Szmrecsanyi. 2007. Recent changes in the function and
13 frequency of standard English genitive constructions: A multivariate analysis of
14 tagged corpora. *English Language and Linguistics* 11: 437–474.
- 15 Hosmer, David and Stanley Lemeshow. 2000. *Applied Logistic Regression*. New York:
16 Wiley.
- 17 Hothorn, Torsten, Kurt Hornik and Achim Zeileis. 2006. Unbiased Recursive
18 Partitioning: A Conditional Inference Framework. *Journal of Computational and*
19 *Graphical Statistics* 15(3): 651–674.
- 20 Hundt, Marianne. 2004. The passival and the progressive passive: a case study in
21 layering in the English aspect and voice systems. In Hans Lindquist and Christian
22 Mair, eds. *Corpus Approaches to Grammaticalization in English*. Amsterdam:
23 Benjamins, 79–120.
- 24 Hundt, Marianne. 2007. *English Mediopassive Constructions*. Amsterdam: Rodopi.
- 25 Hundt, Marianne. 2013. Using web-based data for the study of global English. In
26 Manfred Krug and Julia Schlüter, eds. *Research Methods in Language Variation and*
27 *Change*. Cambridge: University Press, 158–177.
- 28 Hundt, Marianne and Christian Mair. 1999. Agile and uptight genres: The corpus-based
29 approach to language change in progress. *International Journal of Corpus*
30 *Linguistics* 4(2): 221–242.
- 31 Hundt, Marianne, Lena Zipp and André Huber. 2015. Attitudes towards varieties of
32 English in Fiji: A shift to endonormativity? *World Englishes* 34(3): 688–707.
- 33 Hundt, Marianne, Gerold Schneider and Elena Seoane. 2016. The use of the *be*-passive
34 in academic Englishes: Local vs. global usage in an international language. *Corpora*
35 11(1): 31–63.
- 36 Kachru, Yamuna. 2006. *Hindi*. Amsterdam: Benjamins.
- 37 Keenan, Edward L. 1985. Passive in the world's languages. In Timothy Shopen, ed.
38 *Language Typology and Syntactic Description*, Volume 1. Cambridge: Cambridge
39 University Press, 243–281.
- 40 Kondo, Takako. 2005. Overpassivization in second language acquisition. *International*
41 *Review of Applied Linguistics in Language Teaching (IRAL)* 43: 129–161.
- 42 Kortmann, Bernd and Benedikt Szmrecsanyi. 2009. World Englishes between
43 simplification and complexification. In Thomas Hoffmann and Lucia Siebers, eds.
44 *World Englishes. Problems, Properties and Prospects*. Amsterdam: Benjamins, 263–
45 286.
- 46 Leech, Geoffrey, Marianne Hundt, Christian Mair and Nicholas Smith. 2009. *Change in*
47 *Contemporary English: A Grammatical Study*. Cambridge: Cambridge University

- 1 Press.
- 2 Lynch, John. 1998. *Pacific Languages: An Introduction*. Honolulu: University of Hawai'i
- 3 Press.
- 4 Maratsos, Michael. 1988. Crosslinguistic analysis, universals, and language
- 5 acquisition. In Frank E. Kessel, ed. *The Development of Language and Language*
- 6 *Researchers. Essays in Honor of Roger Brown*. Hilldale: Lawrence Erlbaum, 121–
- 7 152.
- 8 Matthews, Stephen and Virginia Yip. 1994. *Cantonese: A Comprehensive Grammar*.
- 9 London and New York: Routledge.
- 10 McFarland, Curtis D. 2008. Linguistic diversity and English in the Philippines. In Ma.
- 11 Lourdes S. Bautista and Kingsley Bolton, eds. *Philippine English: Language and*
- 12 *Literary Perspectives*. Hong Kong: Hong Kong University Press, 131–156.
- 13 Menard, Scott. 2010. *Logistic Regression: From Introductory to Advanced Concepts*
- 14 *and Applications*. Thousand Oakes: SAGE Publications.
- 15 Pinheiro, José C. and Douglas M. Bates. 2000. *Mixed-effects models in S and S-PLUS*.
- 16 New York: Springer.
- 17 R Core Team . 2016. *R: A language and environment for statistical computing*. R
- 18 Foundation for Statistical Computing, Vienna, Austria.
- 19 <https://www.R-project.org/>.
- 20 Ransom, Evelyn. 1979. Definiteness and animacy constraints on passive and double-
- 21 object constructions in English. *Glossa* 13: 215–240.
- 22 Rosenbach, Anette. 2007. Animacy and grammatical variation: findings from English
- 23 genitive variation. *Lingua* 118: 151–171.
- 24 Sandahl, Stella. 2000. *A Hindi Reference Grammar*. Leuven: Peeters.
- 25 Schachter, Paul and Fe T. Otanes. 1972. *Tagalog Reference Grammar*. Berkeley:
- 26 University of California Press.
- 27 Schneider, Edgar W. 2007. *Postcolonial English: Varieties around the World*.
- 28 Cambridge: Cambridge University Press.
- 29 Schneider, Gerold. 2008. *Hybrid Long-Distance Functional Dependency Parsing*. PhD,
- 30 University of Zurich.
- 31 Schütz, Albert J. 2014. *Fijian Reference Grammar*. Honolulu: Pacific Voices Press.
- 32 Seoane, Elena. 2006. Changing styles: on the recent evolution of scientific British and
- 33 American English. In Christiane Dalton-Puffer, Dieter Kastovsky, Nikolaus Ritt and
- 34 Herbert Schendl, eds. *Syntax, Style and Grammatical Norms: English from 1500-*
- 35 *2000*. Bern: Peter Lang, 191–211.
- 36 Seoane, Elena. 2009. Syntactic complexity, discourse status and animacy as
- 37 determinants of grammatical variation in English". *English Language and*
- 38 *Linguistics* 13(3): 365–384.
- 39 Seoane, Elena. 2012. Givenness and word order: A study of long passives in Modern
- 40 and Present-Day English. In Anneli Meurman-Solin, María José López-Couso and
- 41 Bettelou Los, eds. *Information Structure and Syntactic Change in the History of*
- 42 *English*. Oxford: Oxford University Press, 139–163.
- 43 Seoane, Elena and Lucía Loureiro-Porto. 2005. On the colloquialization of scientific
- 44 British and American English. *ESP Across Cultures* 2: 106–118.
- 45 Shibatani, Masayoshi. 1988. Voice in Philippine languages. In Masayoshi Shibatani, ed.
- 46 *Passive and Voice*. Amsterdam: Benjamins, 85–142.
- 47 Silverstein, Michael. 1976. Hierarchy of features and ergativity. In Robert M. W. Dixon,

- 1 ed. *Grammatical Categories in Australian Languages*. Canberra: Australian Institute
2 of Aboriginal Studies, 112–171.
- 3 Strobl, Carolin, Torsten Hothorn & Achim Zeileis. 2009. Party on! A new, conditional
4 variable-important measure for random forests available in the party package. *The*
5 *R Journal* 1(2): 14–17.
- 6 Szmrecsanyi, Benedikt, Jason Grafmiller, Benedikt Heller and Melanie Röthlisberger.
7 2016. Around the world in three alternations: modelling syntactic variation in
8 varieties of English. *English World-Wide* 37(2): 109–137.
- 9 Tagliamonte, Sali & Harald Baayen. 2012. Models, forests, and trees of York English:
10 *Was/were* variation as a case study for statistical practice. *Language Variation and*
11 *Change* 24, 135–178.
- 12 Thomason, Sarah G. 2013. Innovation and contact: the role of adults (and children). In
13 Daniel Schreier and Marianne Hundt, eds. *English as a Contact Language*.
14 Cambridge: Cambridge University Press, 283–297.
- 15 Wasow, Thomas. 2002. *Postverbal Behavior*. Stanford: CSLI Publications.
- 16 Wasow, Thomas and Jennifer Arnold. 2003. Post-verbal constituent ordering in
17 English. In Gunter Rohdenburg and Britta Mondorf, eds. *Determinants of*
18 *Grammatical Variation in English*. Berlin: de Gruyter, 119–154.
- 19 Xiao, Richard, Tony McEnery and Yufang Qian. 2006. Passive constructions in English
20 and Chinese. A corpus-based contrastive study. *Languages in Contrast* 6(1): 109–
21 149.
- 22 Zaenen, Annie, Joan Bresnan, M. Catherine O’Connor, Jean Carletta, Andrew Koontz-
23 Garboden, Tom Wasow, Gregory Garretson and Tatiana Nikitina. 2004. Animacy
24 encoding in English: Why and how. In *Proceedings of the ACL-04 Workshop on*
25 *Discourse Annotation*.
26 (<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.154.7>; accessed 23 November
27 2017)
- 28 Zipp, Lena. 2014. *Educated Fiji English. Lexico-grammar and Variety Status*.
29 Amsterdam: Benjamins.
- 30 Zuur, Alain F., Elena Ieno, Neil J. Walker, Anatoly A. Saveliev, and Graham M. Smith.
31 2009. *Mixed Effects Models and Extensions in Ecology with R*. New York: Springer.