

# About text frequencies in historical linguistics\*

Benedikt Szmrecsanyi  
Freiburg Institute for Advanced Studies

April 2, 2012

This paper is concerned with the limitations of inferring grammar change from variable text frequencies in historical corpus data. We argue that fluctuating frequencies of grammatical variants in real time are a function not only of changing grammars, but are also conditioned by what we call ‘environmental’ changes (for example, content changes) that affect the textual habitat. As a case study, we explore the English genitive alternation in the Late Modern English period and demonstrate that the English *s*-genitive is and always has been preferably used with animate possessors; if for some reason animate NPs are rare in some specific historical period or text, this will trivially depress *s*-genitive rates and boost *of*-genitive rates. Against this backdrop, the paper advocates usage of statistical modeling to probe the probabilistic underpinning of grammatical variability in diachrony, for the sake of keeping apart trivial habitat-induced frequency change and grammar change proper.

## 1. Introduction

This study addresses the problematic status of text frequencies as a diagnostic marker of grammar change in corpus-based, variationist research designs. Any baby – even an otherwise very happy one – will cry when it gets cold. This does not mean, of course, that the baby has changed; what has changed is the baby’s environment, or habitat. We submit that corpus-based historical linguists often face a similar issue, in that fluctuating frequencies of grammatical variants are a function not only of changing grammars, but are also conditioned by environmental changes in the textual habitat. So the crucial problem is that diachronically variable text frequencies often entangle environmental differences and grammatical changes. To disentangle the two types of change, we will argue that instead of focussing solely on text frequencies (how *often* do language users use some

---

\* I am grateful for the feedback to an earlier version of this paper presented at the 2011 Boston Workshop on ‘How can new corpus-based techniques advance historical description and linguistic theory?’. The usual disclaimers apply. This material is based upon work supported by the National Science Foundation under Grant No. BCS-1025602.

linguistic variant?), analysts need to marshal statistical modeling to explore the possibly historically evolving probabilistic conditioning of variants (*why* do language users use the variants that they use?). This approach yields a more reliable diagnostic of grammar change.

Our proposal is not actually a very original one – variationist sociolinguists, for example, have been long aware that language change may manifest in extremely subtle shifts in the stochastic effects of conditioning factors, and that mere text frequencies (or variant rates, for that matter) may be as much about culture as they are about linguistics. Yet in the corpus-based historical linguistics community, we are dealing with a deeply entrenched reliance on the diagnostic power of corpus frequencies. We hasten to add that all other things being equal, diachronically fluctuating text frequencies of grammatical forms may indeed point to grammar change – but especially in historical linguistics, all other things are rarely equal. This is the central point that the present study seeks to emphasize.

This paper is structured as follows: In Section 2, we conduct a thought experiment to illustrate the problem. Section 3 further sets the stage by offering some crucial definitions and by presenting a very concise review of the relevant literature. Section 4 discusses variability between the *s*-genitive and the *of*-genitive in the Late Modern English period as a case study that highlights the limited potential of text frequencies as a diagnostic of grammar change. Section 5 offers some concluding remarks.

## 2. A thought experiment

Let us assume an ancient culture with some language  $X$  that has (at least in principle) an overt future marker  $F$ . The historical record that survives consists exclusively of one particular text type  $R$ . Assume further that before time  $t_1$ , a religious norm outlaws talk about the future in this particular text type. At time  $t_1$ , though, this norm is rescinded – the baby gets cold. The result is a beautiful *s*-curve-shaped frequency explosion of future marker  $F$ . Consider Figure 1, which plots hypothetical text frequencies of  $F$  on the  $y$ -axis<sup>1</sup> against hypothetical units of time on the  $x$  axis. Prima facie, the curve in Figure 1 looks like a language change phenomenon of the sort that every historical linguist is dreaming of discovering once in her lifetime. But we know how this curve came about, and so we must wonder: is the frequency change depicted in Figure 1 a symptom of grammar change, or should we be careful and avoid the term ‘grammar change’ in this context? This question is a rhetorical one, of course. Few analysts would regard the curve in Figure 1 as having anything to do with grammar change – rather, it depicts the time course of an environmental change that has altered the subject matters in one particular text type. It is clear that in the long run the frequency change depicted in Figure 1 may very well have linguistic and/or grammatical consequences (see, e.g., Bybee 2006), but in the short run the curve in Figure 1 is what it is: trivial in linguistic terms.

---

1 Note that for our thought experiment and for the argument in this paper, it does not matter at all if we consider absolute frequencies (e.g. frequency per million words of running text) or some relative frequency measure (e.g. the rate of  $F$  usage vis-à-vis usage of some other grammatical marker, such as the present tense).

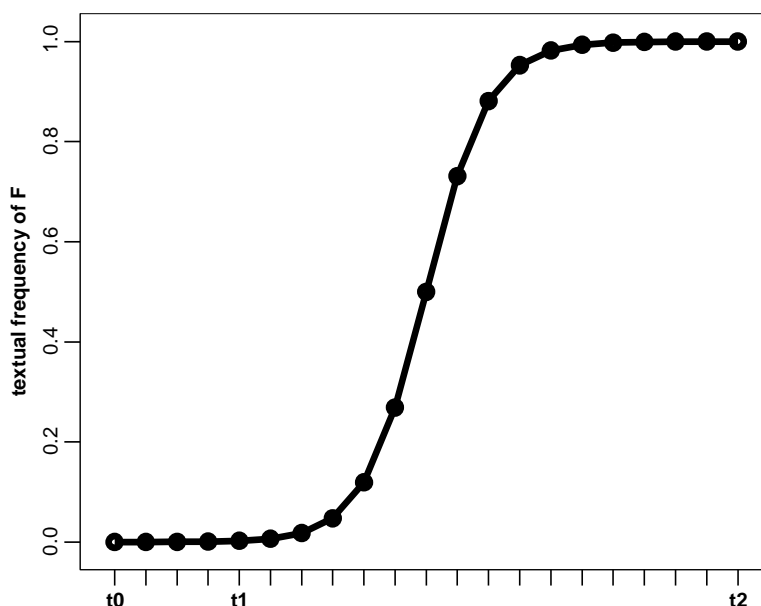


Figure 1: Hypothetical frequencies of future marker  $F$  ( $y$ -axis) at times  $t_0$ ,  $t_1$ , and  $t_2$  ( $x$ -axis).

### 3. Preliminaries: Definitions and literature review

The present study falls within the remit of variationist (socio)linguistics in a modern and fairly broad sense (Labov, 1982; Tagliamonte, 2001; Bresnan and Ford, 2010). This means that we follow sociolinguistic theory and recent probabilistic approaches to language in assuming that grammatical variation and change are often, and maybe even typically, gradient rather than categorical in nature. But what is ‘grammar’? We define grammar as the knowledge of

1. a set of *basic structural units*, however one’s theoretical framework may conceive of them – symbols (e.g. verbs), constructions (e.g. the *s*-genitive construction), or linguistic variables à la Weiner and Labov (1983) (e.g. the active-passive alternation);
2. a set of *probabilistic constraints*, also known as variable rules (in the spirit of Sankoff and Labov 1979; Bresnan and Hay 2008). Probabilistic constraints may read as follows: *an inanimate genitive p’or decreases the odds that speakers of English will use the inflectional genitive by 95%*.

Note that probabilistic constraints may come with  $p = 1$  (i.e. certainty), in which case we are not dealing with choice processes but with categorical constraints fully compatible with categorical notions of grammars in the spirit of e.g. Chomsky (1965) (for example, *in subordinate clauses, speakers of German use verb-final word order with probability 1*).

So, against this backdrop, what is grammar *change*? Trivially, a grammar change may change the set of basic structural units; so, for example, lexical units may develop into new grammatical units. This sort of change is extremely well studied (e.g. Hopper and Traugott 1993). But we will assume that grammar change may also be probabilistic in nature, altering the *constraint set-up* – and possibly *ranking* – that fuels linguistic choice-making. Crucially, probabilistic change ought to be *habitat-independent*, in that it should be a general change observable across text types. This is another way of saying that the frequency profile of grammatical variants generated by the alleged change should not be tied to, or predictable on the basis of, particularities of a specific text or genre. We also believe that real grammatical change should operate *below the level of conscious awareness* (parlance of Labov 1972). In regard to the habitat-independence and non-awareness that we are requiring, Sapir’s notion of ‘drift’ seems apropos:

The drift of a language is constituted by the *unconscious* selection on the part of its speakers of [...] variations [...] In the long run any new feature of the drift becomes part and parcel of the *common, accepted speech*. (Sapir 2004: 127; emphases mine)

Given these criteria, we stress that the frequency change in our thought experiment (Figure 1) clearly does not qualify as grammar change. First, the frequency explosion in our thought experiment is tied to one particular text type in which for non-linguistic reasons, usage of the future marker was non-existent before time  $t_1$ . Second, the booming usage of the future marker after time  $t_1$  is clearly not engendered by “drifty” linguistic choices below the level of consciousness. Instead, our hypothetical writers consciously chose to write about new subject matters (that is, the future); this was an environmental change that secondarily triggered increased text frequencies of future marker  $F$ .

To what extent is the foregoing discussion relevant to current studies in corpus-based historical linguistics? We begin by acknowledging the primacy of speech: the most appropriate text type to study grammar change is, or would be, everyday face-to-face interaction, a genre that drives change but is simultaneously fairly immune to stylistic fads and fashions (for example, Paul 1920: 32). In our thought experiment, we may suppose that speakers of language  $X$  presumably had always used future marker  $F$  in ordinary conversation (remember that it was text type  $R$  only which was subject to religious censorship). The problem, of course, is that the long-term historical record that we have in historical linguistics at the present time does not document speech. In the absence of historical records of face-to-face interaction, however, the dynamics of written registers (read: environmental changes affecting the textual habitat) are a serious confounding variable – much as in our thought experiment, the norm outlawing talk about the future is a confounding factor. Again, we stress that this problem is in principle well-known (see, e.g., Biber and Finegan 1989; Biber and Conrad 2001; Hundt and Mair 1999). Yet, there has been a tendency to shrug off the problem, and/or to concede defeat by accepting that trying to disentangle grammar change from environmental change (such as dynamically evolving written text types) is a hopeless endeavor:

the history of written English is the history of registers [...] all this variation

[...] [is] part of the history of a language – not as “a language” in abstraction but as a language as it is used. (Curzan, 2009: 1103)

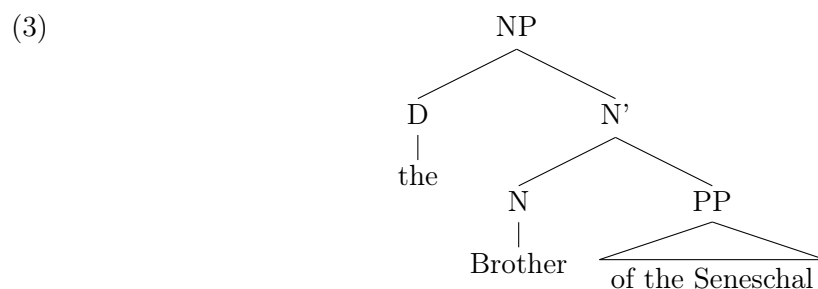
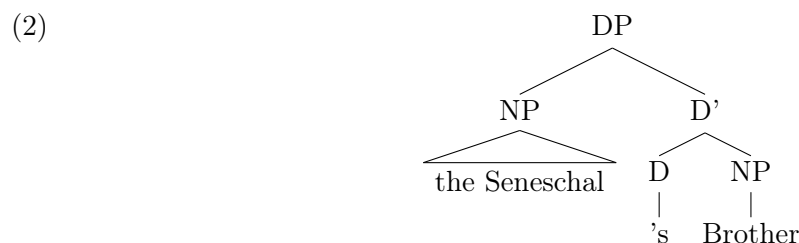
The present study is an extended empirical argument that we can do better than that, and that it is in fact possible to disentangle environmental change (for example, fads and fashions affecting a written register) from grammar change proper. To do so, though, we need to go beyond mere discussions of text frequencies.

#### 4. A case study: genitive variability in Late Modern English

To furnish a little case study, we discuss in what follows variability between the *s*-genitive (as in (1-a)) and the *of*-genitive (as in (1-b)) in the Late Modern English period:

- (1) a. before **[the Seneschal]s [Brother]** could arrive, he was secured by the Governor of Newport <1682pro1.n2b><sup>2</sup>  
 b. the Duke of Norfolk, having lately received another Challenge from **[the Brother]** of **[the Seneschal]**, went to the place appointed <1682pro1.n2b>

The genitive alternation is a primarily syntactic alternation where the order of the so-called possessor (*the Seneschal*) and the so-called possessum (*(the) brother*) can be switched around. In (2) and (3) we report two contemporary syntactic analyses (Radford 1988, 1990), which demonstrate that the two genitive variants are clearly different beasts syntactically: the *s*-genitive is analyzed as a DP (where possessive 's is a determiner and the head of the entire phrase) whereas the *of*-genitive is a straightforward NP with a PP nested into it.



<sup>2</sup> All linguistic examples in this paper are drawn from the ARCHER corpus (see Section 4.2) and are referenced by ARCHER text identifiers.

#### 4.1. The history of genitive variability

Historically, the *of*-genitive is the incoming form, which appeared during the ninth century. According to Thomas (1931: 284) (cited in Mustanoja 1960: 75), the inflected genitive vastly outnumbered the periphrasis with *of* up until the twelfth century. In the Middle English period, we begin to observe “a strong tendency to replace the inflectional genitive by periphrastic constructions, above all by periphrasis with the preposition *of*” (Mustanoja 1960: 70), to the extent that the inflected genitive came close to extinction (Jucker 1993: 121). The frequencies calculated by Thomas (1931) show that by the fourteenth century, the *of*-genitive had a market share of about 84%. Yet in the Early Modern English period, we see a revival of the *s*-genitive, “against all odds” (Rosenbach 2002: 184). In Present-Day English, empirical research has reported comparatively high frequencies of the *s*-genitive (for example, Rosenbach 2002; Szmrecsanyi and Hinrichs 2008). The consensus is that the *s*-genitive is spreading right now (for example, Potter 1969; Dahl 1971; Raab-Fischer 1995; Rosenbach 2003; Szmrecsanyi 2009, 2010).

#### 4.2. Data

We re-analyze the genitive dataset presented in Wolk et al. (submitted) (which in turn partially overlaps with the dataset investigated in Szmrecsanyi in press). The dataset is drawn from ARCHER, *A Representative Corpus of Historical English Registers*, release 3.1 (Biber et al. 1994). ARCHER covers the period between 1650 and 1999, spans about 1.8 million words of running text, and samples a number of different registers. We shall restrict attention to *s*-genitives and *of*-genitives in ARCHER’s British English news (a fairly ‘agile’ genre [Hundt and Mair 1999: 236]) and letters section (a register that is considered fairly oral, at least in regard to private letters [Raumolin-Brunberg 2005: 57]). Our textual database thus comprises 257 texts and totals roughly 242,000 words of running text spread out fairly evenly over the real time periods sampled in ARCHER. Note that the corpus comes subdivided into seven a priori periods of 50 years.

#### 4.3. The linguistic variable

Recall that this study is situated in the variationist sociolinguistics framework. This is why right at the outset, we need to circumscribe the variable context to define interchangeable genitive contexts in which either genitive construction is acceptable (Labov 1966, 1972). Using \*’s, *of*, and \*s as search strings, Wolk et al. (submitted) manually extracted, in a strictly semasiological fashion, all occurrences matching the following patterns:

- [full NP]’s [full NP without determiner];
- [full NP]s [full NP without determiner];
- [full NP]’ [full NP without determiner];
- [full NP] *of* [full NP].

Subsequently, Wolk et al. used a detailed coding scheme to eliminate non-interchangeable genitive contexts (e.g. *of*-constructions with modifying function, or *of* tokens that are part of titles [e.g. *the king of England*]). In this endeavor, Wolk et al. established criteria on the basis of previous research on genitive variation (e.g. Rosenbach 2002; Hinrichs and Szmrecsanyi 2007). These criteria yield genitive constructions that are interchangeable in principle, rather than necessitating a coder's intuition. The end product is a dataset consisting of  $N = 3,824$  interchangeable genitives covering the period between 1650 and 1999.

#### 4.4. Annotation

Next, the 3,824 genitive observations in the dataset were richly annotated for a range of contextual constraints (or: conditioning factors), including the following:

**Possessor animacy.** According to the literature, this is the most crucial constraint on the genitive alternation: the more human and animate a possessor, the more likely it is to take the *s*-genitive (for example, Altenberg 1982; Rosenbach 2008). The annotation distinguishes between five handcoded possessor animacy categories (coding scheme: Zaenen et al. 2004):

1. animate possessors (e.g. *the Bishop's personal security squad* <1979stm2.n8b>)
2. collective possessors (e.g. *the Gentlemen of the Academy* <1723dai2.n3b>)
3. time possessors (e.g. *yesterday's outbreaks* <1967stm1.n8b>)
4. locative possessors (e.g. *the inhabitants of this island* <1872gla1.n6b>)
5. inanimate possessors (e.g. *the rays of greatness* <1748ches.x3b>)

**Genitive relation.** Genitive constructions can encode a variety of semantic relations. Wolk et al. (submitted) follow Rosenbach (2002) and distinguish prototypical and non-prototypical genitive relations. Prototypical relations – which according to the literature favor the *s*-genitive – include:

1. legal ownership relations (e.g. *Mr Ian Smith's cattle ranch and farm* <1979stm1.n8b>)
2. body part relations (e.g. *the murderers legs* <1653merc.n2b>)
3. kinship relations (e.g. *the Duke of Berwick's Son* <1715eve1.n3b>)
4. part-whole relations (e.g. *The Hull of a Ship* <1735rea1.n3b>)

**Constituent length.** According to the so-called 'principle of end-weight', in VO languages such as English speakers and writers tend to place longer, heavier constituents after shorter, lighter ones (for example, Behaghel 1909; Grafmiller and Shih 2011). Thus if the possessor is heavy, there should be a general preference for the *of*-genitive because it places the possessor last; if the possessum is heavy, a general preference

	<i>of-genitive</i>		<i>s-genitive</i>		Total		corpus (words)	size
1650-1699	312	(69%)	139	(31%)	451	(100%)		35k
1700-1749	364	(71%)	152	(29%)	516	(100%)		34k
1750-1799	418	(79%)	109	(21%)	527	(100%)		35k
1800-1849	558	(89%)	70	(11%)	628	(100%)		35k
1850-1899	446	(80%)	109	(20%)	555	(100%)		34k
1900-1949	435	(76%)	134	(24%)	569	(100%)		34k
1950-1999	357	(62%)	221	(38%)	578	(100%)		34k
Total	2890	(76%)	934	(24%)	3,824	(100%)		242k

Table 1: Interchangeable genitive frequencies – absolute (not normalized) versus relative (rates, in %) – by ARCHER period (figures from Wolk et al. submitted).

for the *s-genitive* is expected. Wolk et al. (submitted) operationalized constituent length as a constituent’s length in graphemic characters.

**Final sibilancy.** A final sibilant in the possessor NP (as in *the preparation of this despatch*, <1833tim2.n5b>) discourages usage of the *s-genitive* due to a haplogy or horror aequi effect (for example, Zwicky 1987; Shih et al. to appear). Wolk et al. (submitted) annotated all genitives in the dataset as to whether the possessor phrase ends in a sibilant.

The Wolk et al. paper describes the annotation procedure and the coding schemes on which it is based in exquisit detail. Suffice it to note here that all of the above constraints are non-deterministic in nature – that is, they tend to favor or disfavor particular genitive outcomes in a probabilistic, not categorical, fashion.

#### 4.5. A frequency discussion

In this section, we canvas the frequency distribution of *s-genitives* and *of-genitives* in the Late Modern English period (1650-1999). Table 1 reports raw frequencies (not normalized) as well as genitives rates per ARCHER period (recall that each of these periods covers 50 years). Thus the textual database for the 1650-1699 period (first row) spans about 35,000 words of running text; in this subcorpus, we find in all 451 genitive constructions. 312 of these are *of-genitives*, and 139 are *s-genitives*. In terms of relative frequencies, we are therefore talking about an *of-genitive* rate of 69% and an *s-genitive* rate of 31% in the first ARCHER period.

From reading the literature on long-term historical genitive variability – what with the *s-genitive*’s comeback during the Early Modern English period and its popularity in Present-Day English (see Section 4.1) – one could have expected to see a gradual linear expansion of *s-genitive* rates during the Late Modern English period. Observe now that no such linear expansion emerges from Table 1. What we find instead is a V-shaped pattern:

The *s*-genitive started out with a share of 31% in the 1650-1699 period. Frequencies then started to decline substantially in the 1750-1799 period and hit rock bottom in the 1800-1849 period (11%). Subsequently, *s*-genitive rates recovered such that with a market share of 38% in the 1950-1999 period, the *s*-genitive is more popular now than ever. We note that the V-shaped pattern emerges from relative genitive rates as well as absolute genitive frequencies. Also, the *s*-genitive collapse in the nineteenth century is unlikely to be a sampling issue, as the total number of observations in ARCHER's middle periods is not any lower than in the other periods.

The curious V-shaped frequency pattern that is on display in Table 1 is clearly in need of explanation. Why was the *s*-genitive so unpopular in the nineteenth century? Let us begin by exploring an account that relies on text frequencies as a reliable diagnostic of grammar change (see Szmrecsanyi in press for an in-depth discussion). We premise that the *s*-genitive is often discussed as a counterexample to the unidirectionality (less grammatical > more grammatical) of grammaticalization, having developed from a well-behaved inflection in Old English times to a more clitic-like marker in Present-Day English. Hence, the history of the English *s*-genitive has been cited as an example of “degrammaticalization” (cf., e.g., Janda 1980; Newmeyer 1998) or even “antigrammaticalization” (Haspelmath 2004). Now, the workhorse diagnostic in the corpus-based grammaticalization literature is a construction's overall text frequency (for example, Krug 2000; Mair 2004): “[l]ack of paradigmatic variability [...] accounts for the ubiquity of a feature in the texts of a language” (Lehmann 1995: 142). This is why (at least so the argument goes) “sheer textual frequency is prima facie evidence of degree of grammaticalization” (Hopper and Traugott 1993: 110). In this view, we would diagnose degrammaticalization of the *s*-genitive between 1650 and approximately 1850 (which is when *s*-genitive rates collapsed), and grammaticalization after 1850, and especially during the 20<sup>th</sup> century (which is when *s*-genitive rates recovered substantially).

What is wrong with this account? For starters, the back and forth between degrammaticalization and grammaticalization in a period of merely 350 years is a bit odd (but then again, strange things do happen). The more severe problem is that the frequencies shown in Table 1 demonstrably entangle language-internal developments (for example, (de-)grammaticalization and such like) and language-external developments. Recall (Section 4.4) that animacy of the possessor NP is one of the most crucial conditioning factors in the genitive alternation: animate and/or human NPs favor the *s*-genitive strongly, while inanimate NPs favor the *of*-genitive. And the fact of the matter is that we observe substantial, environmentally induced variability in the distribution of animacy categories (both in terms of genitive NPs and in terms of NPs in general) in written texts during the Late Modern English period. Here is the evidence: Figure 2 presents two area plots that depict the market share of five animacy categories (*y*-axis, in %) against real time (*x*-axis) in ARCHER's news section. The left plot restricts attention to genitive possessor NPs; the right plot is based on a random sample drawn from the general population of NPs (i.e. not necessarily genitive NPs) in ARCHER news. In both samples, animate nouns are on the decline while non-animate nouns are on the rise. In plain English: in the news genre, topics have shifted – from the discussion of animate

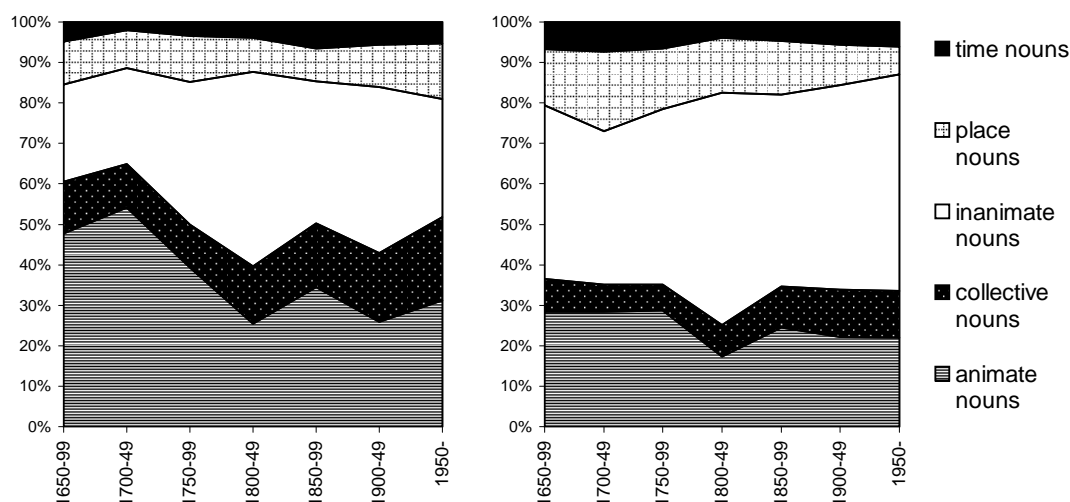


Figure 2: Distribution of animacy categories ( $y$ -axis, in %) against real time ( $x$ -axis) in ARCHER news. Left: genitive possessors only. Right: general noun population in ARCHER news (based on a random sample [ $N = 5, 174$ ]).

entities (as in (4-a)) to the discussion of non-animate entities, such as collective bodies (as in (4-b))

- (4)
- a. They daily expect here Plenipotentiaries from holland, with a final answer upon **the Kings last Propositions** <1672lon2.n2b>.
  - b. Mina had resigned his command, and **the orders of the Executive Government** were duly obeyed by the local authorities and people of Corunna. <1822eva1.n5b>

Hence given that the distribution of animacy categories, which boil down to the prime predictor of genitive choice, is unstable in the textual habitat, it should surprise no one that genitive frequencies fluctuate to some extent. The question is: how much of the real-time variance in genitive frequencies can we explain away by considering this habitat instability? To address this issue, we fit a very simple binary logistic regression model (Pampel 2000) that seeks to predict each of the 3,824 genitive outcomes in the dataset *solely on the basis of the animacy status of the possessor*, following the five-fold categorization (animate/collective/time/place/inanimate) presented in Section 4.4.<sup>3</sup> By considering possessor animacy – and possessor animacy only – we aim to regress out unstable animacy distributions from the frequency picture. Given its simplicity, the resulting model has a surprisingly good fit (Somers  $D_{xy} = 0.64$ ) and captures about 40% of the variability in the dataset (Nagelkerke  $R^2 = 0.386$ ). Next, we take the model's 3,824 genitive outcome predictions and calculate from these predictions mean predicted  $s$ -genitive rates for each of the seven ARCHER periods. These we plot against observed

<sup>3</sup> We utilize R (R Development Core Team 2010) package `lrm`, library `Design`.

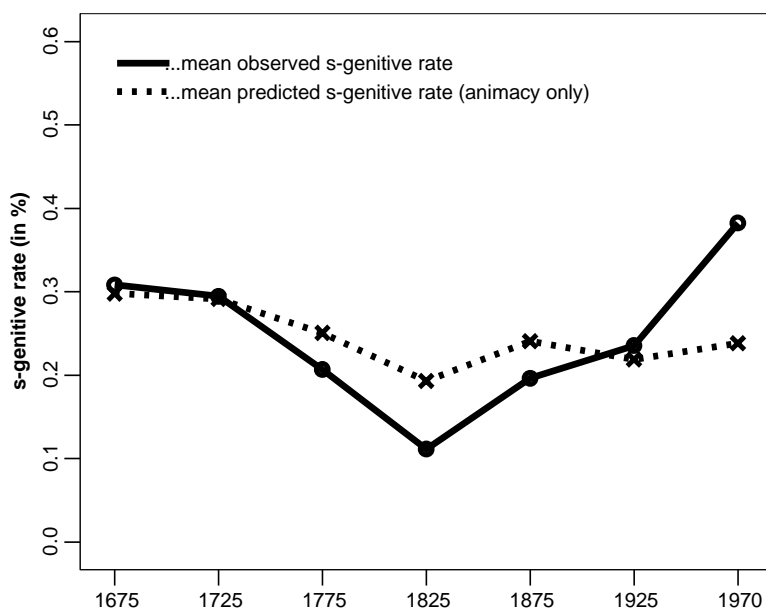


Figure 3: Mean rate of the *s*-genitive (in %, vis-à-vis the *of*-genitive, on *y*-axis) against real time (*x*-axis) (level of granularity: 7 ARCHER periods). Heavy line: observed rates. Dotted line: rates predicted by an animacy-only logistic regression model.

*s*-genitive rates in Figure 3.

Figure 3 shows that the animacy-only model (dotted line) goes some way towards predicting the actually observed collapse of *s*-genitive rates in the nineteenth century (heavy line). Although the animacy-only model admittedly somewhat underestimates the extent of the collapse, it does to some extent account for the frequency decline of the *s*-genitive between 1650 and 1850. This is another way of saying that this part of the story has not much to do with grammar change; what happened – plain and simple – is that a fairly stable grammar of genitive choice produced fewer *s*-genitives because writers chose to write less about animate NPs than they used to. Coming back to the definitions offered in Section 3, we are dealing with a habitat-dependent change that has neither changed the set of basic structural units, nor did it necessarily alter the constraint set-up fueling linguistic variability.

That said, we also note that the animacy-only model (that is to say: the habitat-dependence account) completely fails to predict the comeback of the *s*-genitive during the twentieth century. So something else must have happened in addition to the habitat change, and the task before us is to explore what this other development may have been.

## 4.6. Advanced statistical modeling: mixed-effects binary logistic regression analysis

We have seen in the precious section that unstable animacy distributions in the textual habitat explain away a good deal of the diachronic variability in genitive outcomes; to judge from the animacy-only model's Nagelkerke  $R^2$  value, about 40% of the frequency variability is epiphenomenal. But is the remainder of the frequency variability trivial too? To check if the phenomenal comeback of the *s*-genitive during the twentieth century is a genuine grammar change phenomenon, we will now move on to fit a considerably more sophisticated *mixed-effects logistic regression* model (Pinheiro and Bates, 2000).<sup>4</sup> The model seeks to predict genitive outcomes in the dataset on the basis of

1. a variety of language-internal constraints (also known as *fixed effects*), such as the conditioning factors (possessor animacy, genitive relation, constituent length, and final sibilancy) discussed in Section 4.4 above;<sup>5</sup>
2. interactions between the language-internal constraints and the language-external variable real time (modeled as a scalar predictor);
3. non-repeatable random effects, which control for nuisance factors such as author idiosyncracies and possessor lemma effects.

Wolk et al. (submitted) describe the fixed and random effect structure of the model, as well as the model fitting procedure, in ample detail. Suffice it here to spell out our crucial assumption, in line with the definition offered in Section 3: we can posit (probabilistic) grammar change if the stochastic effect of language-internal predictor variables varies as a function of real time. In other words: if we find that the effect of a language-internal constraint is temporally unstable in a statistically significant way, we can legitimately diagnose grammar change (see Gries and Hilpert 2010 for a similar approach). Note that this is a very elegant and precise criterion which injects some, we believe, welcome methodological rigor into the study of grammatical change.

The resulting model is a very accurate one: it correctly predicts 91.9% of all genitive outcomes in the dataset, and comes with an excellent Somers  $D_{xy}$  value of 0.93. Figure 4 depicts a probabilistic blueprint of genitive choice in ARCHER. The Figure restricts attention to the four conditioning factors discussed in Section 4.4 and reports so-called *odds ratios* (ORs), which quantify the magnitude and the direction of the effect of each predictor on genitive outcomes. ORs specifically indicate how the presence or absence of a feature (for categorical conditioning factors) or how a one-unit increase in a scalar conditioning factor influences the odds for an outcome. Because odds ratios can take values between 0 and infinity, three cases can be distinguished: (i) if  $OR < 1$ , the conditioning factor makes a specific outcome less likely; (ii) if  $OR = 1$ , the conditioning factor has no effect whatsoever on the outcome; (iii) if  $OR > 1$ , the conditioning factor makes a specific outcome more likely (notice that the outcome predicted in Figure 4 is the *of*-genitive). So we observe that all non-animate possessor classes make the

<sup>4</sup> To fit this model, we utilized the R package `lme4`.

<sup>5</sup> In addition, the model considers definiteness of the possessor NP; see Wolk et al. (submitted) for details.

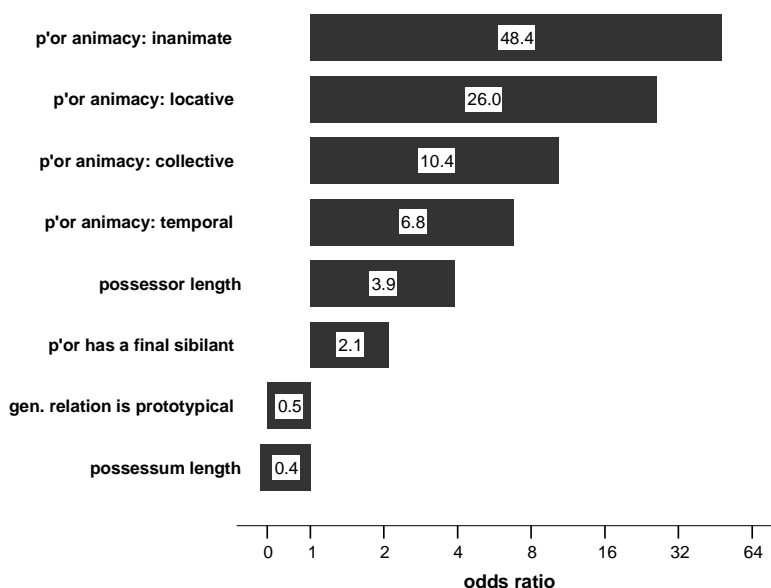


Figure 4: Main effects of genitive predictors: odds ratios (ORs) in logistic regression. Predicted odds are for the *of*-genitive. ORs > 1: favoring; OR < 1: disfavoring. Default levels: Animacy – possessor is human; final sibilancy – possessor does not have a final sibilant; genitive relation – non-prototypical (model parameters adapted from Wolk et al. submitted).

*of*-genitive more likely; for example, if the possessor is inanimate (e.g. *rock*) instead of animate (e.g. *Tom*), the odds for the *of*-genitive increase by a factor of 48.4. Increasingly long possessors and the presence of a final sibilant in the possessor NP also make the *of*-genitive more likely. By contrast, prototypical genitive relations (e.g. kinship – *Tom's brother*) disfavor the *of*-genitive and hence favor the *s*-genitive: if the genitive relation is prototypical instead of non-prototypical, the odds for the *of*-genitive decrease by 50%. Also, increasingly long possessums favor the *s*-genitive.

These findings are as such nothing to write home about – all of the effect directions are the theoretically expected ones, given the literature. But what about diachronic change? Note now that the blueprint in Figure 4 is an estimate for the year AD 1800. Upon closer inspection, it turns out that the probabilistic effect of some of the conditioning factors depicted interacts significantly with real time. In short: we are indeed dealing with genuine (probabilistic) grammar change, according to the criteria that guide this study (Section 3). Among other things, we obtain a significant interaction between real time and possessor animacy. The technicalities need not concern us here (see Wolk et al. submitted for a discussion), but what has happened in a nutshell is that starting in the second half of the nineteenth century, the *s*-genitive has been becoming less strongly disfavored with collective, locative and temporal possessors. This is another way of saying that the grammatical animacy constraint has been weakened. For example, whereas in the

year 1800, a locative possessor (e.g. *the inhabitants of this island*) increased the odds for an *of*-genitive by a factor of 26 (see Figure 4), the corresponding factor in the year 1999 – two centuries later – is only 6.5. Make no mistake: this odds ratio still robustly favors the *of*-genitive. But the motto in probabilistic grammar is that many a little makes a mickle (especially since collective and temporal possessors have also come to disfavor the *s*-genitive less forcefully), and so a subtle change in the probabilistic constraint set-up fueling genitive variation has engendered a robust frequency change. The result is a substantial frequency boost of the *s*-genitive after 1850.

## 5. Conclusion

We have seen that the story of genitive frequencies in the Late Modern English period is complicated. A good deal of the diachronic frequency variability in the dataset can be traced back to environmental changes in the textual habitat (and does not, therefore, diagnose grammar change); but the remainder of the frequency variability is in large part due to (probabilistic) grammar change proper. What happened is that *s*-genitive rates collapsed between 1650 and 1850 because animate NPs became rarer in the textual habitat. So this is the part of the story that has nothing to do with grammar change – the baby just got cold. However, *s*-genitive frequencies recovered between 1850 and 1999 because the grammatical animacy constraint was weakened, such that it became increasingly acceptable to use the *s*-genitive with non-animate possessors. This is the part of the story that indeed involves (probabilistic) grammar change – the baby changed.

So the upshot is that frequency shifts do not always reliably diagnose grammar change. To reiterate, text frequencies of the *s*-genitive collapsed in the nineteenth century because news writers, in particular, wrote less and less about animate entities. This naturally depressed the frequency of the *s*-genitive, a construction which used to be very unpopular with non-animate possessors. It of course remains true that in the long run, environmental changes affecting the textual habitat (such as news writers' increasing interest for non-animate NPs) may very well percolate into grammar (along the lines of Bybee 2006). For example, one may speculate that the frequency collapse of animate possessor NPs in the textual habitat actually triggered the subsequent relaxation of the animacy constraint. But our goal in this paper was not to identify the ultimate causes of grammar change; instead, we sought to establish the extent to which we can trust in text frequencies as a diagnostic marker of grammar change. And in this spirit our case study showed that we can and should differentiate between habitat-induced, environmentally-induced frequency fluctuation and grammar change-induced frequency fluctuation. To accomplish this differentiation – which advances historical description and linguistic theory – we suggested to go beyond a mere discussion of text frequencies, and to explore instead the probabilistic conditioning of grammatical variability. If language users change the way in which they choose variants, then – and only then – can we explain fluctuating text frequencies as the outcome of grammar change.

So the verdict is that text frequencies are a regrettably unreliable and inconclusive

diagnostic of grammar change: they are at best inconclusive, and at worst misleading. The reason is that the all-other-things-being-equal condition, which typically underpins frequency-driven reasoning about grammar change, is rarely met in historical linguistics (and not only in historical linguistics – see Levshina et al. to appear for similar problems in cross-variety analyses). Thus, we advocate conservatism; before positing grammar change, the analyst needs to rule out alternative explanations. In this context, Occam’s razor comes in handy: the principle reminds us to choose the simplest explanation consistent with the facts – and grammar change, alas, is typically *not* the simplest explanation of frequency fluctuation.

## References

- Altenberg, B. (1982). *The Genitive v. the Of-Construction. A Study of Syntactic Variation in 17th Century English*. Malmö: CWK Gleerup.
- Behaghel, O. (1909). Beziehungen zwischen Umfang und Reihenfolge von Satzgliedern. *Indogermanische Forschungen* 25(110-142).
- Biber, D. and S. Conrad (2001). Register variation: A corpus approach. In D. Schiffrin, D. Tannen, and H. Hamilton (Eds.), *The handbook of discourse analysis*, pp. 175–196. Oxford: Blackwell.
- Biber, D. and E. Finegan (1989). Drift and the evolution of English style: a history of three genres. *Language* 65(3), 487–517.
- Biber, D., E. Finegan, and D. Atkinson (1994). ARCHER and its challenges: compiling and exploring A Representative Corpus of Historical English Registers. In U. Fries, G. Tottie, and P. Schneider (Eds.), *Creating and using English language corpora: papers from the Fourteenth International Conference on English Language Research and Computerized Corpora*, pp. 1–13. Amsterdam: Rodopi.
- Bresnan, J. and M. Ford (2010). Predicting syntax: Processing dative constructions in American and Australian varieties of English. *Language* 86(1), 186–213.
- Bresnan, J. and J. Hay (2008). Gradient grammar: An effect of animacy on the syntax of give in New Zealand and American English. *Lingua* 118(2), 245–259.
- Bybee, J. L. (2006). From usage to grammar: The mind’s response to repetition. *Language* 82, 711–733.
- Chomsky, N. (1965). *Aspects of the theory of syntax*. Cambridge, MA: MIT Press.
- Curzan, A. (2009). Historical corpus linguistics and evidence of language change. In A. Lüdeling and M. Kytö (Eds.), *Corpus Linguistics. An International Handbook*, Handbücher zur Sprache und Kommunikationswissenschaft/Handbooks of Linguistics and Communication Science, pp. 1091–1109. Berlin, New York: Mouton de Gruyter.
- Dahl, L. (1971). The s-genitive with non-personal nouns in modern English journalistic style. *Neuphilologische Mitteilungen* 72, 140–172.
- Grafmiller, J. and S. Shih (2011). Weighing in on end weight. *Talk given at the LSA 2011 Annual Meeting, 6-9 January 2011, Pittsburgh, Pennsylvania*.

- Gries, S. T. and M. Hilpert (2010). Modeling diachronic change in the third person singular: a multifactorial, verb- and author-specific exploratory approach. *English Language and Linguistics* 14, 293–320.
- Haspelmath, M. (2004). On directionality in language change with particular reference to grammaticalization. In O. Fischer, M. Norde, and H. Perridon (Eds.), *The nature of grammaticalization*, Typological Studies in Language, pp. 17–44. Amsterdam: Benjamins.
- Hinrichs, L. and B. Szmrecsanyi (2007). Recent changes in the function and frequency of Standard English genitive constructions: a multivariate analysis of tagged corpora. *English Language and Linguistics* 11(3), 437–474.
- Hopper, P. J. and E. C. Traugott (1993). *Grammaticalization*. Cambridge, New York: Cambridge University Press.
- Hundt, M. and C. Mair (1999). 'Agile' and 'uptight' genres: the corpus-based approach to language change in progress. *International Journal of Corpus Linguistics* 4, 221–242.
- Janda, R. D. (1980). On the decline of declensional systems: the overall loss of OE nominal case inflections and the ME reanalysis of -es as his. In E. C. Traugott, R. Labrum, and S. Shepherd (Eds.), *Papers from the 4th International Conference on Historical Linguistics*, pp. 243–252. Amsterdam/Philadelphia: Benjamins.
- Jucker, A. (1993). The genitive versus the of-construction in newspaper language. In A. Jucker (Ed.), *The Noun Phrase in English. Its Structure and Variability*, pp. 121–136. Heidelberg: Carl Winter.
- Krug, M. G. (2000). *Emerging English modals: a corpus-based study of grammaticalization*. Berlin, New York: Mouton de Gruyter.
- Labov, W. (1966). The linguistic variable as a structural unit. *Washington Linguistics Review* 3, 4–22.
- Labov, W. (1972). *Sociolinguistic patterns*. Philadelphia: University of Pennsylvania Press.
- Labov, W. (1982). Building on empirical foundations. In W. Lehmann and Y. Malkiel (Eds.), *Perspectives on Historical Linguistics*, pp. 17–92. Amsterdam, Philadelphia: Benjamins.
- Lehmann, C. (1995). *Thoughts on Grammaticalization*. LINCOS Studies in Theoretical Linguistics. München, Newcastle: LINCOS EUROPA.
- Levshina, N., D. Geeraerts, and D. Speelman (to appear). Towards a 3D-Grammar: variation of Dutch causative constructions. *Journal of Pragmatics*.
- Mair, C. (2004). Corpus linguistics and grammaticalisation theory: statistics, frequencies, and beyond. In C. Mair and H. Lindquist (Eds.), *Corpus Approaches to Grammaticalisation in English*, pp. 121–150. Amsterdam/Philadelphia: Benjamins.
- Mustanoja, T. F. (1960). *A Middle English syntax*, Volume I. Helsinki: Société Néophilologique.
- Newmeyer, F. J. (1998). *Language form and language function*. Cambridge: MIT Press.
- Pampel, F. (2000). *Logistic Regression. A Primer*. Quantitative Applications in the Social Sciences. Thousand Oaks: Sage Publications.

- Paul, H. (1920). *Prinzipien der Sprachgeschichte* (5th ed.). Halle: Niemeyer.
- Pinheiro, J. C. and D. M. Bates (2000). *Mixed-effects models in S and S-PLUS*. New York: Springer.
- Potter, S. (1969). *Changing English*. London: AndrÄf Deutsch.
- R Development Core Team (2010). *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. ISBN 3-900051-07-0.
- Raab-Fischer, R. (1995). Löst der Genitiv die of-Phrase ab? Eine korpusgestützte Studie zum Sprachwandel im heutigen Englisch. *Zeitschrift für Anglistik und Amerikanistik* 43(2), 123–132.
- Radford, A. (1988). *Transformational grammar: a first course*. Cambridge, New York: Cambridge University Press.
- Radford, A. (1990). *Syntactic theory and the acquisition of English syntax: the nature of early child grammars of English* (1. publ. ed.). Oxford: Blackwell.
- Raumolin-Brunberg, H. (2005). The diffusion of subject you: A case study in historical sociolinguistics. *Language Variation and Change* 17, 55–73.
- Rosenbach, A. (2002). *Genitive variation in English: conceptual factors in synchronic and diachronic studies*. Berlin, New York: Mouton de Gruyter.
- Rosenbach, A. (2003). Aspects of iconicity and economy in the choice between the s-genitive and the of-genitive in English. In G. Rohdenburg and B. Mondorf (Eds.), *Determinants Of Grammatical Variation in English*, pp. 379–412. Berlin, New York: Mouton de Gruyter.
- Rosenbach, A. (2008). Animacy and grammatical variation – findings from English genitive variation. *Lingua* 118(2), 151–171.
- Sankoff, D. and W. Labov (1979). On the use of variable rules. *Language in Society* 8, 189–222.
- Sapir, E. (2004). *Language, an introduction to the study of speech*. Mineola, New York: Dover.
- Shih, S., J. Grafmiller, R. Futrell, and J. Bresnan (to appear). Rhythm’s role in genitive construction choice in spoken English. In *Rhythm in Phonetics, Grammar and Cognition*. Berlin, New York: de Gruyter Mouton.
- Szmrecsanyi, B. (2009). Typological parameters of intralingual variability: Grammatical analyticity versus syntheticity in varieties of English. *Language Variation and Change* 21(03), 319–353.
- Szmrecsanyi, B. (2010). The English genitive alternation in a cognitive sociolinguistics perspective. In D. Geeraerts, G. Kristiansen, and Y. Peirsman (Eds.), *Advances in Cognitive Sociolinguistics*, Volume 45, pp. 139–166. Berlin, New York: De Gruyter Mouton.
- Szmrecsanyi, B. (in press). The great regression: genitive variability in Late Modern English news texts. In K. Börjars, D. Denison, and A. Scott (Eds.), *Morphosyntactic categories and the expression of possession*. Amsterdam, Philadelphia: Benjamins.

- Szmrecsanyi, B. and L. Hinrichs (2008). Probabilistic determinants of genitive variation in spoken and written English: a multivariate comparison across time, space, and genres. In T. Nevalainen, I. Taavitsainen, P. Pahta, and M. Korhonen (Eds.), *The Dynamics of Linguistic Variation: Corpus Evidence on English Past and Present*, pp. 291–309. Amsterdam: Benjamins.
- Tagliamonte, S. (2001). Comparative sociolinguistics. In J. Chambers, P. Trudgill, and N. Schilling-Estes (Eds.), *Handbook of Language Variation and Change*, pp. 729–763. Malden and Oxford: Blackwell.
- Thomas, R. (1931). *Syntactical processes involved in the development of the adnominal periphrastic genitive in the English language*. PhD thesis, University of Michigan.
- Weiner, J. and W. Labov (1983). Constraints on the agentless passive. *Journal of Linguistics* 19, 29–58.
- Wolk, C., J. Bresnan, A. Rosenbach, and B. Szmrecsanyi (submitted). Dative and genitive variability in Late Modern English: Exploring cross-constructive variation and change.
- Zaenen, A., J. Carlette, G. Garretson, J. Bresnan, A. Koontz-Garboden, T. Nikitina, M. C. O'Connor, and T. Wasow (2004). Animacy encoding in English: why and how. In D. Byron and B. Webber (Eds.), *Proceedings of the 2004 ACL workshop on discourse annotation, Barcelona, July 2004*, pp. 118–125.
- Zwicky, A. M. (1987). Suppressing the *zs*. *Journal of Linguistics* 23, 133–148.