# An analytic-synthetic spiral
# in the history of English

Benedikt Szmrecsanyi

KU Leuven

Drawing on techniques familiar from quantitative morphological typology (Greenberg 1960), this contribution marshals usage- and frequency-based, aggregative measures of grammatical analyticity and syntheticity to profile the history of grammatical marking in English between circa AD 1100 and AD 1900, tapping into the Penn Parsed Corpora of Historical English series. Results indicate that the post-Old English period is clearly not characterized by a linear drift towards more analyticity and less syntheticity. Instead, analyticity was on the rise until the end of the Early Modern English period, but declined subsequently; the reverse is true for syntheticity. In terms of typological analyticity-syntheticity coordinates, 20th century English texts are actually fairly similar to 12th and 13th century English texts. I suggest that this historical pattern can be interpreted in terms of a Gabelentz (1891)-type spiral.

## 1. Introduction

We teach beginning students of English Language and Linguistics that the history of English is characterized by a drift from synthetic to analytic. The textbook story goes as follows: English has changed

> from a synthetic or inflectional language, which relies on morphological endings to mark grammatical function, to an analytic one, which relies on word order to convey grammatical elations. From a structural point of view this is the most significant change that has occurred in the history of English.    (Fennell 2001: 6)

It is clear that this story is true if we compare Old English to Present-Day English. But as we shall see, when we restrict attention to the post-Old English history of the language, it turns out that there is no longer a linear drift but a cyclical merry-go-round.

This study conceptualizes the analytic-synthetic distinction in a way that makes possible precise and holistic measurements, drawing on quantitative, frequency- and usage-based measures inspired by work in quantitative morphological typology (in particular, Greenberg 1960). I will re-analyze and re-interpret the dataset originally discussed in Szmrecsanyi (2012), where these measure were applied to the Penn Parsed

Corpora of Historical English series. This corpus suite covers the period between circa AD 1100 and AD 1900. The present study thus taps into corpus material to determine the frequency of analytic and synthetic marking, and subsequently calculates an aggregate index of overt grammatical analyticity, which basically measures the text frequency of free grammatical markers, and an aggregate index of overt grammatical syntheticity, which measures the text frequency of bound grammatical markers. So the name of the game in this contribution is calculation of aggregate indices on the basis of naturalistic usage data, inspired by techniques in quantitative morphological typology.

We will see that analyticity peaked towards the end of the Early Modern English period, in the 16th and 17th centuries, and has been on the decline ever since. Conversely, syntheticity had its low point in the Early Modern English period and has been on the increase subsequently. In other words, we seem to be dealing with a cycle (Jespersen 1917) or spiral (Gabelentz 1891) of sorts. One is here, of course, in particular reminded of seminal work by Hodge (1970), who demonstrated that the history of Egyptian is characterized by an analytic > synthetic > analytic cycle.

This article is structured as follows. Section 2 fixes terminology. Section 3 describes the data source. Section 4 discusses the method. Section 5 plots cyclical analyticity-syntheticity fluctuations in a bird's eye perspective. Section 6 explores the linguistic sources of this cyclicity in a jeweler's eye perspective. Section 7 offers some concluding remarks and investigates the extent to which we are dealing with a Gabelentz-Jespersen-Hodge-style cyclical phenomenon.

## 2.    Terminology

The terms "analytic" and "synthetic" are, of course, staple terms in classical work on the cross-linguistic typology of languages. The labels go back to the 19th century; August Wilhelm von Schlegel is usually credited for coining the opposition. I cannot review the rich history of thought in this area (but see Schwegler 1990 for an excellent literature review). A pertinent problem is that the terms "are used in widely different meanings by different linguists" (Anttila 1989: 315). This is why we need to fix terminology right at the outset.

I will be interested, first, in the overt coding of *grammatical* information. Hence, this study will have nothing to say about lexical analyticity and syntheticity. Second, my definition of analytic/analyticity and synthetic/syntheticity is a strictly *formal* one that broadly follows Danchev's definition:

> Formal analyticity evidently implies that the various meanings (grammatical and/or lexical) of a given language unit are carried by two or more free morphemes, whereas formal syntheticity is normally characterized by the presence of one bound morpheme.                                                          (Danchev 1992, 26)

I thus operationally define

– *formal grammatical analyticity* as covering all coding strategies that convey gram-matical information via free grammatical markers, which in turn are defined as synsemantic (see Marty 1908) word tokens devoid of independent lexical mean-ing; and
– *formal grammatical syntheticity* as covering all those coding strategies where grammatical information is signaled by bound grammatical markers.

As for analyticity, this study equates synsemantic word tokens with function (also known as structure or empty) words, which are in the present study taken to be mem-bers of closed word classes: conjunctions (e.g. *and*, *if*), determiners (e.g. *the*), pronouns (e.g. *he*), prepositions (e.g. *in*), infinitive markers (e.g. *to*), modal verbs (e.g. *can*, *will*), and negators (e.g. *not*). This view of analyticity and of what should count as a function word is fairly customary (see standard reference works such as Bussmann, Trauth & Kazzazi 1996: 22 and 471).

With regard to the definition of syntheticity, bound grammatical markers are taken to comprise verbal, nominal, and adjectival inflectional affixes (e.g. past tense *-ed*, plural *-s*, comparative *-er*, and so on), the genitive clitic (as in *Tom's house*), as well as allomorphies including ablaut phenomena (e.g. past tense *sang*), i-mutation (e.g. plural *men*), and other non-regular yet clearly bound grammatical markers. The morphological analysis in this study thus broadly adopts an item-and-process model (Hockett 1954: 396) in which grammatically marked forms are seen as deriv-ing from simple forms via some sort of process. What does not feature in this notion of syntheticity is the "zero morpheme'" construct postulated in some morphological approaches to deal with paradigmatic contrasts in finite verb forms – recall that this study is interested in the *overt* coding of grammatical information. Note also that contracted elements (*'s* as in *it's*, *to* as in *gotta*) do not count as bound markers in the present study's approach. Instead, a form such as *gotta* would be analyzed as consist-ing of two free grammatical markers, *got* and *to*, which may happen to be contracted to various degrees in speech.

As for portmanteau morphemes, this study calculates syntheticity indices in the following fashion: what is measured, in a given textual sample, is not the num-ber of inflectional morphemes per sample (which is what Greenberg's original gross inflectional index measured), but the number of words in a sample that bear *at least* one bound grammatical marker. This is not a trivial adjustment, because depending on one's analytical framework e.g. the form *walks* (as in *he walks the dog*) could be analyzed as containing two grammatical morphemes, {non-past} and {3rd person singular}. But in this study's approach, the form *walks* contains exactly one (overt) grammatical marker, *-s*, which may or may not have more than one meaning.

## 3.   Data

This study's data source is the Penn Parsed Corpora of Historical English series (see ⟨http://www.ling.upenn.edu/hist-corpora/⟩). This corpus suite has the following sub-components: The Penn-Helsinki Parsed Corpus of Middle English, second edition (PPCME2) (Kroch & Taylor 2000); the Penn-Helsinki Parsed Corpus of Early Modern English (PPCEME) (Kroch, Santorini & Diertani 2004); and the Penn Parsed Corpus of Modern British English (PPCMBE) (Kroch, Santorini & Diertani 2010). The three corpora yield a total of 605 texts which span slightly less than four million words of running British English text. Each of the texts in this database can be assigned not only to particular periods in the history of English (Middle English, Early Modern English, Late Modern English), but also to specific centuries, starting with the early twelfth century and ending with the early twentieth century. Notice that the texts represent a variety of text types, such as letters, sermons, handbook prose, and history writing. The texts in the Penn Parsed Corpora of Historical English series are all part-of-speech annotated and syntax-parsed using roughly the same tagsets and annotation schemes. As I will explain in the next section, what will take center stage in the present study is the corpus suite's part-of-speech annotation.

## 4.   Method

This study adopts the method described in detail in Szmrecsanyi (2009: 321–325), which I recapitulate below. The method is inspired by a seminal paper (Greenberg 1960) entitled "A Quantitative Approach to the Morphological Typology of Language," in which Joseph Greenberg in turn drew inspiration from work by Edward Sapir. Greenberg demonstrated that seemingly abstract typological notions are in fact amenable to precise quantitative measurements by calculating a number of indices on the basis of *naturalistic texts.* It thus bears highlighting that Greenberg (1960) is an early exercise in quantitative usage-based typology. Of course, the sample size used in Greenberg (1960) were coherent texts of merely 100 words; to mitigate the problem of point estimates deriving from such small sample sizes, Stepanov (1995) suggested basing the calculation of indices on corpora that "will include hundreds of texts from all existing genres, sources, historical periods etc. as one large sample" (Stepanov 1995: 144). This is precisely the kind of analysis that the present study will conduct.

　　Greenberg specifically defined (i) an index of synthesis, (ii) an index of agglutination, (iii) a compounding index, (iv) a derivational index, (v) a gross inflectional index, (vi) a prefixial index, (vii) a suffixial index, (viii) an isolational index,

(ix) a pure inflectional index, and (x) a concordial index (Greenberg 1960: 187). The *gross inflectional index*, for example, is defined as the number of inflectional morphemes in the analyst's sample divided by the total number of words in the sample (Greenberg 1960: 186–187).

  While faithful to the core idea, I take the liberty adapt Greenberg's method. Greenberg did not calculate a plain analyticity index, and so Kasevič and Jachontov (1982: 37) (cited in Kempgen & Lehfeldt 2004: 1237) proposed an "index of analyticity", which would relate the number of synsemantic words in a given text to the total number of words in that text (see also Kelemen 1970: 62 for a similar proposal). Heeding this proposal, this study will calculate two indices:

1.  The *analyticity index* (henceforth AI), which is defined as the ratio of the number of free grammatical markers in a sample (F) to the total number of words in the sample (W), normalized to a sample size of 1,000 tokens. Hence: AI = F/W × 1,000.
2.  The *syntheticity index* (henceforth SI), which is defined as the ratio of the number of words in a sample that bear a bound grammatical marker (B) to the total number of words in the sample (W), normalized to a sample size of 1,000 tokens. Hence: SI = B/W × 1,000.

The indices have a lower bound of 0, and an upper bound of 1,000 index points.

  How does this study determine the number of free grammatical markers in a text, and the number of words in a text that bear a bound grammatical marker? I exploit the part-of-speech (POS) annotation in the Penn Parsed Corpora of Historical English series, which annotates each individual word token in the corpus database for its word class; this includes information on whether nouns, verbs, adjectives, and certain pronouns carry inflections (a detailed description of the Penn Parsed Corpora of Historical English series' POS tagset, including exemplification, is available at ⟨http://www.ling.upenn.edu/hist-corpora/annotation/toc-long.htm#pos⟩). Given the definition of analyticity and syntheticity detailed above, POS tags (or rather the tokens annotated with POS tags) were subsequently placed into four categories: (i) purely lexical tags, such as singular nouns, which are uninteresting for present purposes; (ii) synthetic tags (essentially all tokens that, following Vennemann 1982, 330, show affixation or mutation to indicate grammatical information), (iii) analytic tags (function words), and (iv) a small number of simultaneously synthetic and analytic tags, such as inflected auxiliary verbs and reflexive pronouns in their plural form. Tables 1 and 2 report the exact tag/token-to-category matches, categorizing analytic tags into eleven broad component categories and synthetic tags into four broad component categories.

**Table 1.** Eleven broad component categories (as defined through POS tags and/or word tokens) loading on the Analyticity Index.[1]

| Component category | POS tag(s) |
| --- | --- |
| 1. conjunctions, complementizers, prepositions, subordinating conjunctions | CONJ (coordinating conjunctions) <br> C (complementizers) <br> P (prepositions, including subordinating conjunctions) |
| 2. determiners | D (determiners) <br> WD (*wh*-determiners) |
| 3. existential *there* | EX (existential THERE) |
| 4. pronouns | PRO (personal pronoun) <br> MAN (indefinite subject pronoun [ME, MAN]) – only in Middle English) <br> WPRO (wh-pronoun) |
| 5. *more/most/less/least* | QR (quantifier, comparative) <br> QS (quantifier, superlative) |
| 6. infinitive markers | TO (infinitival TO, TIL, and AT) <br> FOR (infinitival FOR) <br> FOR+TO (cliticized FOR+TO) |
| 7. modals | MD (modal verb) <br> MD0 (modal verb, untensed) |
| 8. negation | NEG (negation) |
| 9. auxiliary *be* | B*(BE) + verb |
| 10. auxiliary *do* | D*(DO) + verb |
| 11. auxiliary *have* | H*(HAVE) + verb |

**Table 2.** Four broad component categories (as defined through POS tags and/or word tokens) loading on the Syntheticity Index

| Component category | POS tag(s) |
| --- | --- |
| 1. the *s*-genitive | $ (possessive), except for PRO$ |
| 2. inflected comparative and superlative adjectives | ADJR (adjective, comparative) <br> ADJS (adjective, superlative) |
| 3. plural nouns | NPRS (proper noun, plural) <br> NS (common noun, plural) <br> OTHERS (OTHER, nominal use, plural) |

(*Continued*)

---

**1.** Note that expletive *it* cannot be considered here because it does not have a unique POS tag.

**Table 2.** (Continued)

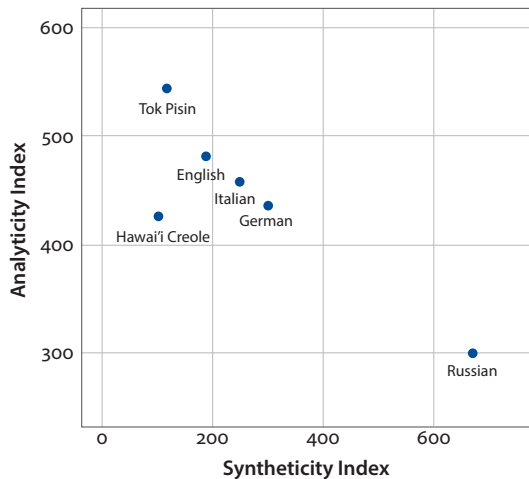| Component category | POS tag(s) |
|---|---|
| 4. inflected lexical and primary verbs | VAG (present participle) |
| | VAN (passive participle, verbal or adjectival) |
| | VBD (past, including past subjunctive) |
| | VBN (perfect participle) |
| | BAG (BE, present participle) |
| | BED (BE, past, including past subjunctive) |
| | BEN (BE, perfect participle) |
| | DAG (DO, present participle) |
| | DAN (DO, passive participle, verbal or adjectival) |
| | DOD (DO, past, including past subjunctive) |
| | DON (DO, perfect participle) |
| | HAG (HAVE, present participle) |
| | HAN (HAVE, passive participle, verbal or adjectival) |
| | HVD (HAVE, past, including past subjunctive) |
| | HVN (HAVE, perfect participle) |
| | VBP (present, including present subjunctive) + inflection |
| | MD (modal verb) + inflection |
| | BEP (BE, present, including present subjunctive) |
| | DOP (DO, present, including present subjunctive) |
| | HVP (HAVE, present, including present subjunctive) |
| | BE (BE, infinitive) ending in *-en* |
| | DO (DO, infinitive) ending in *-en* |
| | HV (HAVE, infinitive) ending in *-en* |
| | VB (infinitive, verbs other than BE, DO, HV) ending in *-en* |



**Figure 1.** Analyticity Index scores (*y*-axis) against Syntheticity Index scores (*x*-axis) in European languages and two English-based creole languages (adapted from Siegel, Szmrecsanyi & Kortmann 2014)

Exemplification can be found in Section 6 below. On the technical plane, a retrieval script written in the programming language Perl automatically established the text frequencies of the relevant POS-tags (or POS-tag categories) in the data set, and calculated the index scores. Subsequently, the quantitative information was analyzed and visualized using the statistical software package R.
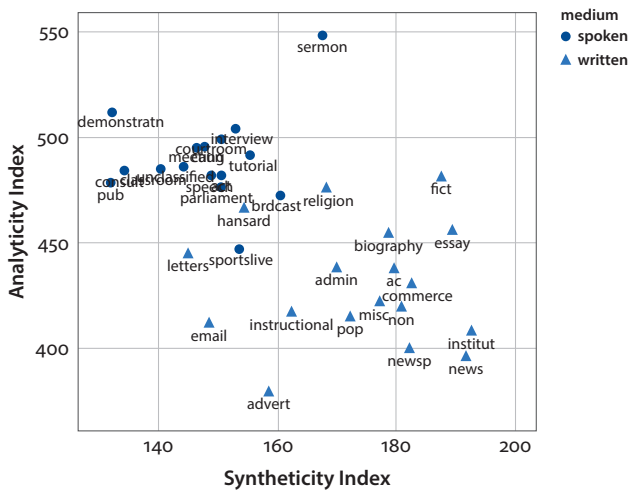


**Figure 2.** Analyticity Index scores (*y*-axis) against Syntheticity Index scores (*x*-axis) of text types sampled in the British National Corpus (adapted from Szmrecsanyi 2009)

To instill confidence in the method, Figures 1 and 2 show how the method rates different languages and text types in contemporary English. Figure 1 locates a number of different European languages (English, Italian, German, and Russian) as well as two English-based creole languages (Tok Pisin and Hawai'I Creole) in a two-dimensional analyticity-syntheticity plane. It turns out that Russian is the most synthetic and least analytic language in the sample, while Tok Pisin is the most analytic and least synthetic language (Hawai'i Creole is also fairly non-synthetic, but less analytic than Tok Pisin). It is probably the case that Figure 1 is a good representation of many linguists' gut feelings about these languages. Figure 2 applies the method to the various spoken and written text types sampled in the British National Corpus (BNC). Here, we find a very clear split between spoken text types and written text types: spoken text types are more analytic and less synthetic than written text types. Again, this is essentially the pattern that most register analysts would expect to see. For more discussion of these patterns, I refer the reader to the papers mentioned in the Figure captions. The crucial point in terms of the present study is that the method seems to work as advertised. With this in mind, I now go on to explore changes, with regard to analyticity and syntheticity, in the history of English.

## 5.    The bird's eye perspective: The big merry-go-round

This section investigates the overall development of the two indices in the eight centuries covered in the Penn Parsed Corpora of Historical English series. Table 3 reports mean Analyticity Index scores and mean Syntheticity Index scores by century, averaging over all text types in the corpus. Figure 3 visually depicts the longitudinal trajectories by plotting index scores ($y$-axis) against real time ($x$-axis) (level of granularity: individual texts), and approximating the relationship by a fit curve.

**Table 3.**  Mean Analyticity Index and Syntheticity Index by century.

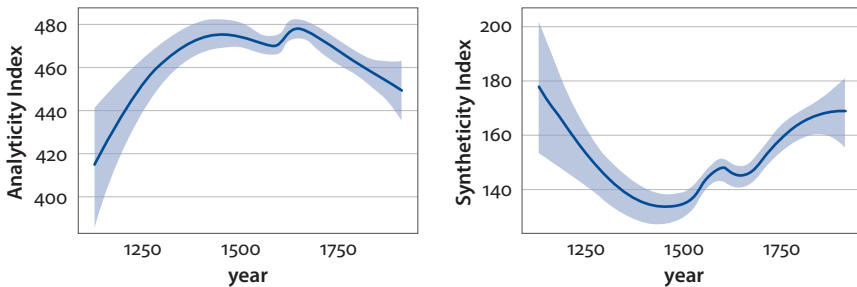| Century | Analyticity Index | Syntheticity Index |
|---|---|---|
| 12th | 449 | 196 |
| 13th | 434 | 151 |
| 14th | 481 | 155 |
| 15th | 470 | 140 |
| 16th | 473 | 141 |
| 17th | 477 | 147 |
| 18th | 464 | 162 |
| 19th | 455 | 166 |
| 20th | 444 | 178 |



**Figure 3.**  Index level variability by year of creation: LOESS smoothed fit curves with confidence region. Database: all texts in the Penn Parsed Corpora of Historical English series. Left: Analyticity Index. Right: Syntheticity Index

Note, first, that analyticity increased between the twelfth and the fourteenth century, remained fairly constant until the seventeenth century, and decreased subsequently. So in the 20th century, we find on average 444 analytic markers per 1,000 words of running text, which is indeed comparable to the levels in the 12th (449) and 13th (434)

centuries. In the 17th century, by contrast, texts in the corpus material exhibit on average no less than 477 analytic markers per 1,000 words of running text.

Syntheticity, in turn, decreased rather robustly between the 12th century, when we find on average 196 inflected words per 1,000 words of running text, and the 15th century, when we find only about 140 inflected words per 1,000 words of running text (see Table 3). However, syntheticity levels rebounded in subsequent centuries. An average 20th century text, for example, in the Penn Parsed Corpus series features no less 178 inflected words per 1,000 words of running text.



**Figure 4.** Mean Analyticity Index scores (*y*-axis) against mean Syntheticity Index scores (*x*-axis) by century. Database: all texts in the Penn Parsed Corpora of Historical English series

Figure 4 is a two-dimensional analyticity-syntheticity plane that visually depicts the cyclical nature of the variability. The diagram highlights the fact that analyticity-syntheticity variability after the Old English period can hardly be described in terms of a steady trend, or drift, towards more analyticity and less syntheticity. It is true that the data point for the 12th century is the most synthetic one in Figure 2. Between the 12th and the 13th century, syntheticity decreases robustly, but so does analyticity. Between the 13th to the 14th century, we see a huge surge in analyticity, while syntheticity levels decrease only slightly. Nothing much then happens between the 14th and the 17th centuries; both analyticity and syntheticity levels remain quite stable. Between the 17th and the 20th centuries, however, we observe a steady and incremental drift towards more syntheticity and less analyticity. The data point for the 20th century is both significantly more synthetic and less analytic than the data point for the seventeenth century.

The upshot is that in terms of analyticity-syntheticity coordinates, the 20th century has come almost full circle back to where we started in the 12th century.
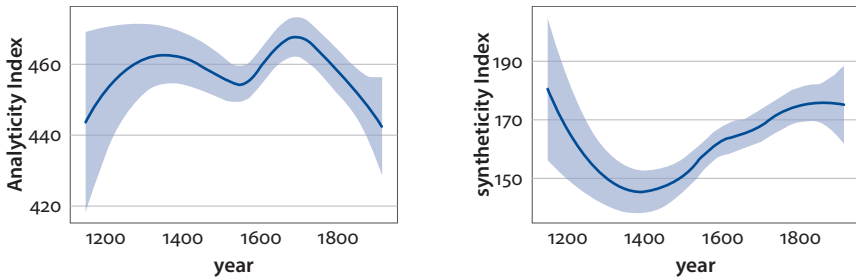
**Figure 5.** Informative texts ($N = 260$) in the Penn Parsed Corpora of Historical English series – index level variability by year of creation: LOESS smoothed fit curves with confidence region. Left: Analyticity Index. Right: Syntheticity Index
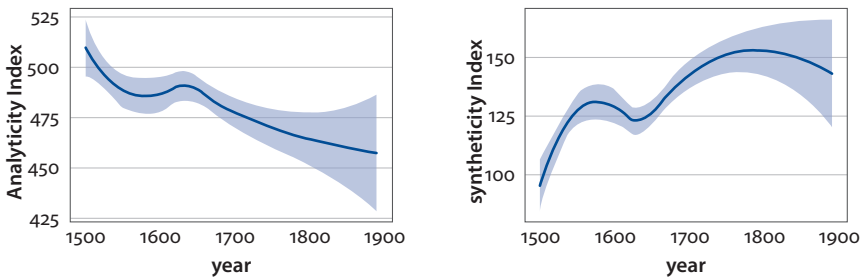
**Figure 6.** Letters ($N = 219$) in the Penn Parsed Corpora of Historical English series – index level variability by year of creation: LOESS smoothed fit curves with confidence region. Left: Analyticity Index. Right: Syntheticity Index
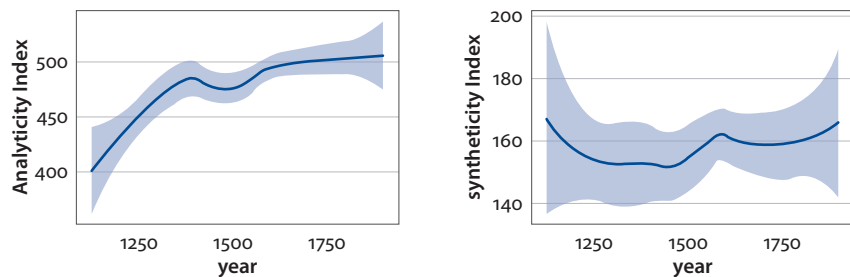
**Figure 7.** Religious texts ($N = 73$) in the Penn Parsed Corpora of Historical English series – index level variability by year of creation: LOESS smoothed fit curves with confidence region. Left: Analyticity Index. Right: Syntheticity Index
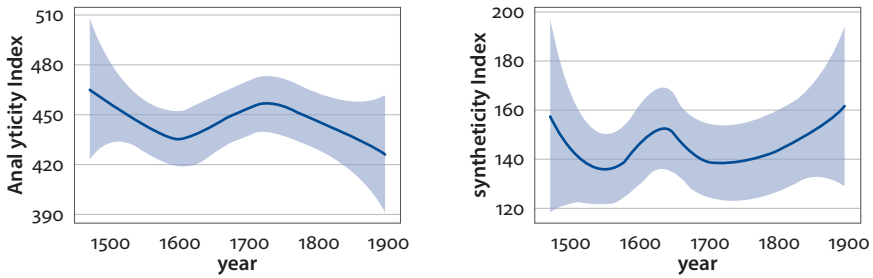
**Figure 8.** Imaginative texts (*N* = 35) in the Penn Parsed Corpora of Historical English series – index level variability by year of creation: LOESS smoothed fit curves with confidence region. Left: Analyticity Index. Right: Syntheticity Index

Is the merry-go-round in Figure 4 somehow an artefact of the design of Penn Parsed Corpora series, which samples a number of different text types? Are we missing something when we aggregate over these text types? To address these particular questions – rather than to characterize different text types in the Penn Parsed Corpora series, which is not my primary concern – Figures 5–8 canvas the development of the indices in the four major text types represented in the Penn Parsed Corpora series: informative texts, letters, religious texts, and imaginative texts. The developments in informative texts (Figure 5) mirror the overall development (see Figure 3), which is not entirely surprising as informative texts constitute text category with the best coverage in the Penn Parsed Corpora series. Letters (Figure 6) show a curious pattern: in this text type analyticity is overall on the decline, while syntheticity is on the rise. Moving on, in both religious texts (Figure 7) and in imaginative texts (Figure 8) we see a weak version of the U-shaped syntheticity pattern familiar from Figure 3. But while analyticity is quite steadily increasing in religious texts, it is fairly stable in imaginative texts, with some ups and downs. In conclusion, the overall trajectories to be found in the Penn Parsed Corpora series (see Figure 3) primarily reflect the developments in one particular text type, informative texts. That said, Figure 9 – which replicates the two-dimensional plane in Figure 2 but excludes informative texts from the calculation – shows that even when informative texts are ignored, there is no linear drift, by any stretch of imagination, from syntheticity to analyticity.

## 6.   The jeweler's eye perspective

The previous section has relied on aggregation to study multi-feature, big-picture patterns. In this section, I trade in the bird's eye perspective for the jeweler's eye perspective, and thus engage in an exercise of index deconstruction: what are the linguistic features that are implicated in the patterns discussed in the previous section?
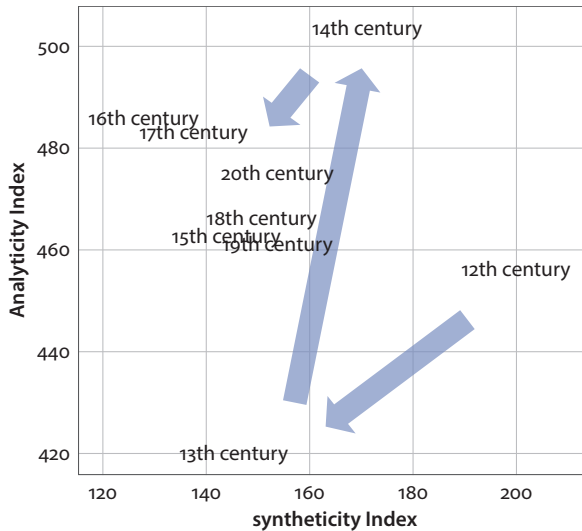
**Figure 9.** Mean analyticity indices (*y*-axis) against mean syntheticity indices (*x*-axis) by century. Database: all texts except informative prose in the Penn Parsed Corpora of Historical English series

In this spirit, Figure 10 plots the frequency trajectories of the 11 features which are loading on the Analyticity Index. The features whose frequency pattern is broadly in line with the inverted U-shaped Analyticity Index trajectory familiar from Figure 3 are the following:

–    Conjunctions, as in (1a), complementizers, as in (1b), prepositions, as in (1c), subordinating conjunctions, as in (1d) (Figure 10a);

   (1a)    He was just getting into Talk with […], **but/CONJ** during the first part of the visit he said very little. (AUSTEN-180X)

   (1b)    The morning was so wet **that/C** I was afraid […] (AUSTEN-180X)

   (1c)    but Frank who alone could go to Church called for her **after/P** service (AUSTEN-180X)

   (1d)    I wished **when/P** I heard them say so, that they could have heard […] (AUSTEN-180X)

–    Pronouns, as in (2) (Figure 10d);

   (2)    They were very civil to **me/PRO**, as they always are (AUSTEN-180X)

–    Infinitive markers, as in (3) (Figure 10f);

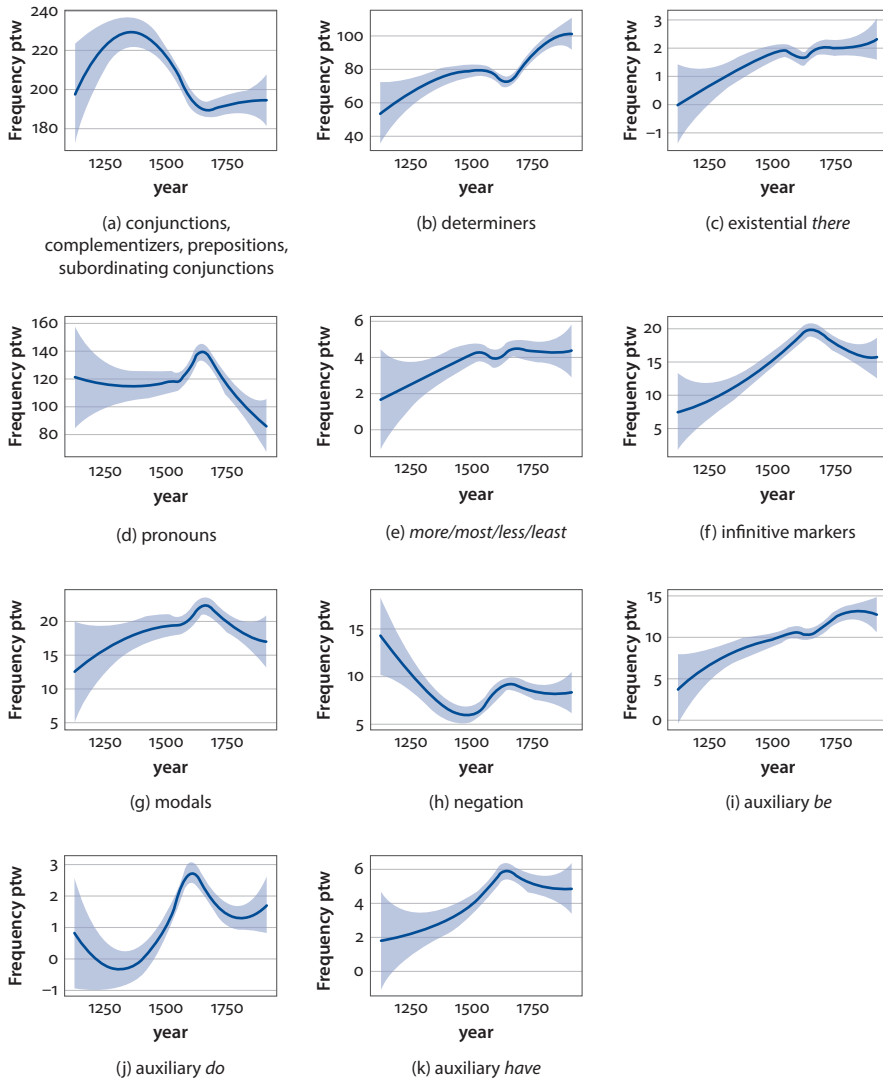   (3)    but this is not thought likely **to/TO** happen (AUSTEN-180X)

**Figure 10.** Frequency per thousand words (ptw) of component features/feature groups loading on the Analyticity Index by year of creation: LOESS smoothed fit curves with confidence region. Database: all texts in the Penn Parsed Corpora of Historical English series

- Modals, as in (4) (Figure 10g);

(4) He **must/MD** think it very strange that I do not acknowledge the receipt […] (AUSTEN-180X)

- Auxiliary *do*, as in (5) (Figure 10j);

(5) My dear Cassandra How **do/DOP** you do ? (AUSTEN-180X)

– Auxiliary *have*, as in    (6) (Figure 10k);

(6)    he & I **have/HVP** practiced together two mornings (AUSTEN-180X)

Negative markers, as in (7), are on the decline (Figure 10h).

(7)    Mrs. E. Leigh did **not/NEG** make the slightest allusion to my Uncle's Business […] (AUSTEN-180X)

By contrast, determiners (Figure 10b and (8)), existential *there* (Figure 10c and (9)), the comparative and superlative markers *more/most/less/least* (Figure 10e and (10)), and auxiliary *be* (Figure 10i and (11)) have all become more frequent over the course of time.

(8)    My Mother wrote to her **a/D** week ago (AUSTEN-180X)

(9)    **There/EX** will then be the Window-Curtains […] (AUSTEN-180X)

(10)    she considers her own going thither as **more/QR** certain […] (AUSTEN-180X)

(11)    A fortnight afterwards she is to **be/BE** called again […] (AUSTEN-180X)

A one-way ANOVA which uses a three-partite sub-corpus distinction (PPCME2 versus PPCEME versus PPCMBE) as grouping variable suggests that the following four analytic features (or feature groups) exhibit the most extensive real-time variance: determiners, pronouns, infinitive markers, and auxiliary verbs.
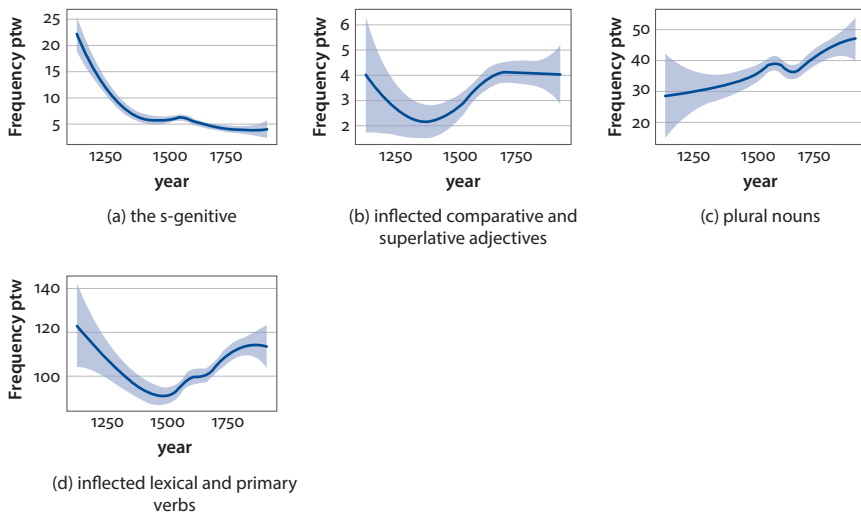


**Figure 11.** Frequency per thousand words (ptw) of component features/feature groups loading on the Syntheticity Index by year of creation: LOESS smoothed fit curves with confidence region. Database: all texts in the Penn Parsed Corpora of Historical English series

Let us turn to the four features that load on the Syntheticity Index; Figure 11 plots the frequencies of these features against real time. Two of these, inflected comparative and superlative adjectives (Figure 11b and (12)) and inflected lexical and primary verbs (Figure 11d and (13)), have U-shaped frequency trajectories that are similar to the overall trajectory of the Syntheticity Index in Figure 3.

(12)   I can no **longer/ADJR** take his part against you, as I did nine years ago (AUSTEN-180X)

(13)   but as she **found/VBD** it agreable I suppose there was no want of Blankets […] (AUSTEN-180X)

The *s*-genitive (Figure 11a and (14)), on the other hand, is sloping downward, while plural nouns (Figure 9c and (15)) are sloping upward.

(14)   The Duke of Gloucester **'s/$** death sets my heart at ease (AUSTEN-180X)

(15)   I shall be very glad to hear from you, that we may know how you all are, especially the two **Edwards/NPRS** (AUSTEN-180X)

## 7.   Discussion and Conclusion

Inspired by techniques developed in quantitative morphological typology (Greenberg 1960), this contribution has sketched the development of grammatical analyticity and syntheticity in the Penn Parsed Corpora of Historical English series, which covers the period between the 12th and the 20th century. The key insights can be summarized as follows. First, the period under study is clearly not characterized by a steady drift towards more analyticity and less syntheticity. Instead, analyticity was on the rise until the end of the Early Modern English period, but declined subsequently; the reverse is true for syntheticity. But with regard to the frequency of analytic versus synthetic marking, 20th century English is quantitatively almost back to the analyticity-syntheticity coordinates that characterize 12th century English. Second, we have seen that this pattern primarily reflects the developments in one particular text type that is extremely well represented in the Penn Parsed Corpora series, informative texts (though I hastened to add that the other text types covered in the corpus material – letters, religious texts, imaginative texts – do not exhibit linear drifts either, as we have seen). Third, a jeweler's eye analysis of the linguistic features that are associated with either index revealed that the overall merry-go-round pattern diagnosed in the bird's eye perspective seems to be a function of the frequency trajectories of the following features: conjunctions, complementizers, prepositions, and subordinating conjunctions; pronouns; infinitive markers; modals; auxiliary *do*; auxiliary have; inflected comparative and superlative adjectives; and inflected lexical and primary verbs.

So the verdict is that when restricting attention to the post-Old English history of the language, we see a good deal of circularity. This circularity no doubt bears a number of resemblances to Gabelentz-Jespersen-Hodge-style cycles or spirals. That said, there are also a number of important differences and caveats that must be considered:

*Time depth*. The dataset I have explored in this contribution covers eight centuries. This is a fairly short interval, compared e.g. to the millennia-encompassing analysis presented in Hodge (1970).

*Frequency changes.* The developments sketched in the present study are primarily about drifts in usage frequency ("fluctuations in analyticity/syntheticity", to borrow a phrase from Schwegler 1990: 191). Additions to, or losses from, the inventory of grammatical markers do not take center stage: eight centuries is simply too short a time span to feature substantial inventory changes.

*Erosion and replacement not mandatory.* It follows that the cyclical developments we are observing in the data are not necessarily about "changes where a phrase or word gradually disappear and is replaced by a new linguistic item" (Gelderen 2009: 2). For example, we saw that nouns bearing plural inflections have become more frequent in the course of the past few centuries. But this does not mean, of course, that the English language has acquired new inflectional markers of nominal plurality. It rather means that for some reason (stylistics, content, discourse pragmatics, etc.), language users have come to use existing plural markers more often. Thus Hodge's motto "one man's morphology was an earlier man's syntax" (Hodge 1970: 3) is a theme that the present contribution's findings do not necessarily speak to. Sources of renewal in the present contribution's approach are possibly old forms that had never died out.

*Aggregation.* To explore big-picture cyclical developments, I eschewed the "single-feature study" (parlance of Nerbonne 2008) perspective advocated by e.g. Heine, Claudi & Hünnemeyer (1991: 246) and instead studied multi-feature cyclical patterns in an aggregate perspective, adopting a whole-language view in the tradition of e.g. Schlegel (1818) and Sapir (1921).

*No discrete steps.* Thanks to the usage-based and frequency-oriented method I have been using, it is difficult to distinguish between discrete steps along the lines of e.g. Hoeksema (2009), who distinguishes four stages in the negative cycle.

*Quantitative versus qualitative change.* Modern analyticity and syntheticity is, of course, qualitatively different from its Early English counterpart. For example, determiners have become increasingly important as an analytic category, but pronouns have been on the decline. Conversely, the possessive marker used to be a more important synthetic marker than it is now, whereas inflected adjectives are on the rise. The point is that contrary developments like these may "gang up", as it were, to

create numerically similar Analyticity and Syntheticity Index scores, although the contribution of particular linguistic features may vary quite dramatically. What we are seeing in the data is thus a spiral (Gabelentz 1891), rather than a cycle.

And this takes us to the most important issue to keep in mind: the indices I have been calculating in this study capture but one characteristic – the typological nature and frequency of grammatical marking – that can be used to compare (historical) language varieties and texts. Take texts cmvices1.m1, from the Penn-Helsinki Parsed Corpus of Middle English, and text benson-1908, from the Penn Parsed Corpus of Modern British English. Both texts exhibit very similar index levels (cmvices1.m1: AI – 458, SI -162; benson-1908: AI – 455, SI – 167), but the texts certainly "feel" extremely different, as a cursory glance at the excerpts in (16) and (17) shows.

(16)   Dies ilche modinesse, +deih hie habbe hlot and dale mang alle o+dre sennes, na+del+as hie haue+d ane, +de is hire swi+de neih and swi+de hersum, +de me haue+d swi+de ofte beswiken, +tat is, Vana Gloria, idel wulder o+der idel +gelp. (cmvices1.m1)

(17)   As regards the art of teaching it is difficult to lay down rules, because every man must find out his own method. It is easy to say that the first requisite is patience, but the statement requires considerable modification. (benson-1908)

And so although 20th century English may be similar, in terms of quantitative analyticity and syntheticity, to 12th and 13th century English, there is no way texts from these periods can be confused.

## Acknowledgments

## References

Anttila, Raimo. 1989. *Historical and Comparative Linguistics* [Current Issues in Linguistic Theory 6]. Amsterdam: John Benjamins. doi:10.1075/cilt.6

Bussmann, Hadumod, Trauth, Gregory & Kazzazi, Kerstin. 1996. *Routledge Dictionary of Language and Linguistics*. London: Routledge.

Danchev, Andrei. 1992. The evidence for analytic and synthetic developments in English. In *History of Englishes: New Methods and Interpretations in Historical Linguistics*, Matti Rissanen, Ossi Ihalainen, Terttu Nevalainen & Irma Taavitsainen (eds), 25–41. Berlin: Mouton de Gruyter.

Fennell, Barbara A. 2001. *A History of English: A Sociolinguistic Approach*. Oxford: Blackwell.

von der Gabelentz, Georg. 1891. *Die Sprachwissenschaft: Ihre Aufgaben, Methoden und Bisherigen Ergebnisse*. Leipzig: Weigel.

van Gelderen, Elly. 2009. Cyclical change. An introduction. In *Cyclical Change* [Linguistik Aktuell/Linguistics Today 146], Elly van Gelderen (ed.), 1–12. Amsterdam: John Benjamins. doi:10.1075/la.146

Greenberg, Joseph H. 1960. A quantitative approach to the morphological typology of language. *International Journal of American Linguistics* 26(3): 178–94.  doi:10.1086/464575

Heine, Bernd, Claudi, Ulrike & Hünnemeyer, Friederike. 1991. *Grammaticalization: A Conceptual Framework*. Chicago IL: University of Chicago Press.

Hockett, Charles F. 1954. Two models of grammatical description. *Word* 10: 210–231.

Hodge, Carleton T. 1970. The linguistic cycle. *Language Sciences* 13: 1–7.

Hoeksema, Jack. 2009. Jespersen recycled. In *Cyclical Change* [Linguistik Aktuell/Linguistics Today 146], Elly van Gelderen (ed.), 15–34. Amsterdam: John Benjamins. doi:10.1075/la.146.04hoe

Jespersen, Otto. 1917. *Negation in English and Other Languages*. Copenhagen: Host.

Kasevič, Vadim & Jachontov, Sergej E. (eds). 1982. *Kvantitativnaja Tipologija Jazykov Azii I Afriki* (A Quantitative Typology of Asian and African Languages). Leningrad.

Kelemen, József. 1970. Sprachtypologie und Sprachstatistik. In *Theoretical Problems of Typology and the Northern Eurasian Languages*, László Dezső & Peter Hajdú (eds), 53–63. Amsterdam: Gruener.

Kempgen, Sebastian & Lehfeldt, Werner. 2004. Quantitative typologie. In *Morphologie. Ein Internationales Handbuch Zur Flexion Und Wortbildung*, Geert E. Booij, 1235–1246. Berlin: Mouton de Gruyter.

Kroch, Anthony, Santorini, Beatrice & Diertani, Ariel. 2004. *Penn-Helsinki Parsed Corpus of Early Modern English*. ⟨http://www.ling.upenn.edu/hist-corpora/PPCEME-RELEASE-2/index.html⟩

Kroch, Anthony, Santorini, Beatrice & Diertani, Ariel. 2010. *Penn Parsed Corpus of Modern British English*. ⟨http://www.ling.upenn.edu/hist-corpora/PPCMBE-RELEASE-1/index.html⟩

Kroch, Anthony & Taylor, Ann. 2000. *Penn-Helsinki Parsed Corpus of Middle English,* 2nd edn. ⟨http://www.ling.upenn.edu/hist-corpora/PPCME2-RELEASE-3/index.html⟩

Marty, Anton. 1908. *Untersuchungen zur Grundlegung der Allgemeinen Grammatik und Sprachphilosophie*. Halle: Niemeyer.

Nerbonne, John. 2008. Variation in the aggregate: An alternative perspective for variationist linguistics. In *Northern Voices: Essays on Old Germanic and Related Topics Offered to Professor Tette Hofstra*, Kees Dekker, Alasdair MacDonald & Hermann Niebaum (eds), 365–82. Leuven: Peeters.

Sapir, Edward. 1921. *Language: An Introduction to the Study of Speech*. New York NY: Harcourt, Brace and Company.

von Schlegel, August Wilhelm. 1818. *Observations Sur La Langue et La Littérature Provençales*. Paris: Librairie grecque-latine-allemande.

Schwegler, Armin. 1990. *Analyticity and Syntheticity: A Diachronic Perspective with Special Reference to Romance Languages*. Berlin: Mouton de Gruyter.  doi:10.1515/9783110872927

Siegel, Jeff, Szmrecsanyi, Benedikt & Kortmann, Bernd. 2014. Measuring analyticity and syntheticity in Creoles. *Journal of Pidgin and Creole Languages* 29(1): 49–85. doi:10.1075/jpcl.29.1.02sie

Stepanov, Arthur V. 1995. Automatic typological analysis of Semitic morphology. *Journal of Quantitative Linguistics* 2(2): 141–150.  doi:10.1080/09296179508590043

Szmrecsanyi, Benedikt. 2009. Typological parameters of intralingual variability: Grammatical analyticity versus syntheticity in varieties of English. *Language Variation and Change* 21(3): 319–53.  doi:10.1017/S0954394509990123

Szmrecsanyi, Benedikt. 2012. Analyticity and syntheticity in the history of English. In *The Oxford Handbook of the History of English*, Terttu Nevalainen & Elisabeth Closs Traugott (eds), 654–665. Oxford: OUP.

Vennemann, Theo. 1982. Isolation – Agglutination – Flexion? Zur Stimmigkeit Typologischer Parameter. Fakten und Theorien. In *Festschrift für Helmut Sinn Zum 65. Geburtstag*, Sieglinde Heinz & Ulrich Wandruszka (eds), 327–34. Tübingen: Narr.