



HAL
open science

Towards an n-grammar of English

Bert Cappelle, Natalia Grabar

► **To cite this version:**

Bert Cappelle, Natalia Grabar. Towards an n-grammar of English. Constructionist Approaches to Second Language Acquisition and Foreign Language Teaching, 2016. hal-01426700

HAL Id: hal-01426700

<https://hal.science/hal-01426700>

Submitted on 4 Jan 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40



IV Constructing a constructicon for L2 learners

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

1 Bert Cappelle and Natalia Grabar

2 **Towards an n-grammar of English**

3
4

5 **Abstract:** In this chapter, it is shown how we can develop a new type of learner's
6 or student's grammar based on n-grams (sequences of 2 or 3, 4, etc. items) auto-
7 matically extracted from a large corpus, such as the Corpus of Contemporary
8 American English (COCA). The notion of n-gram and its primary role in statistical
9 language modelling is first discussed. The part-of-speech (POS) tagging provided
10 for lexical n-grams in COCA is then demonstrated to be useful for the identifica-
11 tion of frequent structural strings in the corpus. We propose using the hundred
12 most frequent POS-based 5-grams as the content around which an 'n-grammar'
13 of English can be constructed. We counter some obvious objections to this
14 approach (e.g. that these patterns only scratch the surface, or that they display
15 much overlap among them) and describe extra features for this grammar, relat-
16 ing to the patterns' productivity, corpus dispersion, functional description and
17 practice potential.

18
19 **Keywords:** ESL/EFL, POS n-grams, frequency, construct-i-con, grammar
20 teaching

21
22

23 **1 Introduction: Words, words, words, but where's** 24 **the grammar?**

25
26

27 Linguists these days are being spoiled with increasingly large corpora. There is
28 for instance Oxford University's popular British National Corpus (BNC), which
29 contains 100 million words and which is freely available from Mark Davies's
30 website, among other online services.¹ Davies's bigger and more up-to-date Corpus
31 of Contemporary American English (COCA) contains 450 million words (Davies
32 2008-) and his more recently added Global Web-Based English (GloWbE) allows
33 us to search through 1.9 billion words (Davies 2013). This web corpus is now
34 dwarfed by others, such as ENCOW14, which contains almost 17 billion tokens.²
35 And then there is the biggest 'corpus' of all, the indexable part of the World
36 Wide Web itself, which as long ago as June 2006 was estimated to contain 14.3
37 billion web pages and to increase in size by 280 million web pages a day (De

38
39 ¹ <http://corpus.byu.edu/>, last accessed on 2 February 2015.

40 ² <http://corporafromtheweb.org/encow14/#more-72>, last accessed on 28 February 2015.

1 Kunder 2006). Whether we use a search engine such as Google or query a
2 comparatively much smaller but still very large corpus designed for linguistic
3 research, what we have at our finger tips in each case is a venerable treasure
4 trove of data about real language use.

5 The availability of frequency-based word lists compiled from such large
6 corpora of varied texts (e.g. Davies and Gardner 2010) may be of great benefit
7 to practitioners in the field of teaching English as a second or foreign language
8 (ESL/EFL). And indeed, for several decades, corpora have already served as a
9 valuable aid in developing vocabulary teaching materials (see, e.g., McCarthy
10 and O'Dell 2001 for a well-known product). Corpus-based vocabulary teaching
11 prevents certain 'pet' expressions in ESL/EFL, such as *raining cats and dogs*,
12 from being taught too vigorously, and common but less favorite ones, such as
13 *right up your (or his, her, etc.) alley*, from being ignored altogether.

14 In sharp contrast to the teaching of lexis, grammar teaching does not
15 involve much attention to frequency and focuses instead on, for example, how
16 to construct interrogative or passive structures from canonical (declarative,
17 active) ones. Very often, though, grammar is not even taught that explicitly,
18 since this is felt to go against the prevailing functionally-oriented approach
19 to language learning. It is our impression that when grammar is taught at all,
20 explicitly or in task-based learning settings, the sequence and selection of
21 grammar patterns is mostly a matter of convention and convenience.

22 Lexis and grammar, as we shall have the opportunity to see, are two sides of
23 the same coin, in that concrete lexical items (words and collocations) belong
24 to more abstract categories (word classes and phrasal structures). One might
25 therefore assume that teaching specific words and expressions automatically
26 results in teaching rules of grammar. Moreover, as there are patterns which
27 combine concrete and more abstract pieces, a distinction between lexis and
28 grammar is often claimed to be illusionary (cf. Ellis and Cadierno 2009). Never-
29 theless, abstract structures also have an existence which is not wholly reducible
30 to the collocations and idioms that they represent. This is because grammar
31 patterns are generalizations not just over idioms but over lexically rather
32 mundane combinations as well. For instance, the passive construction is not
33 'just' used in expressions such as *to be cast in stone* or *to be caught between a*
34 *rock and a hard place*. It is a structure which can be applied productively and it
35 should therefore be taught as such. So, since abstract phrasal constructions do
36 not only underlie frequently used lexical sequences but also provide blue-prints
37 for creative combinations, they need to be focused on in their own right. A
38 purely lexical approach cannot suffice in language teaching.

39 Most importantly, we need to know which abstract structures are most fre-
40 quent in the language, because as it is, ESL/EFL is still in dire need of a reliable,

1 ordered inventory of the most frequently used grammatical patterns in English.
 2 Material developers would much appreciate linguists to provide them with a list
 3 of common grammar structures for active mastery, to be distinguished from less
 4 common patterns that learners can acquire more incidentally. This is, in any
 5 case, what the first author of this paper has heard first-hand from an educational
 6 advisor for Flemish secondary school teachers (Johan Delbaere, personal com-
 7 munication). The result of this lack of an objective standard of frequent patterns
 8 is that constructions that are typically taught may not actually be that frequent
 9 and, conversely, that frequent constructions may go unnoticed by material
 10 developers. The aim of this paper, therefore, is to show that we can exploit
 11 corpus data not just to identify frequent lexical items and lexical patterns of
 12 co-occurrence but also to find frequent grammar patterns. That is, just as lexicog-
 13 raphers have been successful in detecting common words and collocations,
 14 grammarians should really start using corpora to find the most common struc-
 15 tural patterns in a language.

16 To be fair, some existing grammars do take corpus frequencies into account.
 17 A prime example is the *Longman Grammar of Spoken and Written English* (Biber
 18 et al. 1999), which is entirely corpus-based and provides detailed frequency
 19 information (across registers), as well as the *Longman Student Grammar of*
 20 *Spoken and Written English* (Biber, Conrad, and Leech 2002), which is based on
 21 the latter. Another example of a corpus-based grammar is Cobuild's two-volume
 22 *Grammar Patterns* (Francis, Hunston, and Manning 1996, 1998), whose lexico-
 23 grammatical approach, outlined in Hunston and Francis (2000), is heavily influ-
 24 enced by work by Halliday and Sinclair (e.g. Halliday 1978; Sinclair 1991). Other
 25 early studies that comment on frequency of use are referenced in Celce-Murcia
 26 and Larsen-Freeman (1999). However, despite these valuable works, linguists
 27 so far have not yet produced any *ranking* of frequent grammar patterns for the
 28 benefit of EFL/ESL teachers, students and material developers.

29 We will show that this can be achieved by using *n-grams* – continuous
 30 sequences of *n* (i.e. any specified number of) items. Our demonstration will be
 31 restricted to n-grams extracted from the COCA corpus. This is entirely for practical
 32 reasons, as will become clear.³ We believe that common lexical and grammatical
 33 n-grams are constructions, in a Construction Grammar sense: they are form-

35
 36 ³ Apart from COCA, there are other corpora which allow n-gram-based grammar studies. For
 37 instance, as shown in Cappelle (2014), using Google's *Ngram Viewer* (Michel et al. 2010), we
 38 can exploit the n-grams extractable from Google Books for (diachronic) research into grammar
 39 patterns, since this corpus has been tagged and allows part-of-speech searches (Lin et al. 2012).
 40 The COW corpora also provide n-gram data sets (<http://hpsg.fu-berlin.de/cow/ngrams/>, last
 accessed on 28 February 2015).

1 function pairings which native speakers have memorized (and which learners of
 2 a language should acquire) as a result of their high frequency. For Construction
 3 Grammarians, frequency is only one of the criteria to identify constructions,
 4 another possible criterion being the unpredictable nature of the link between a
 5 unit's form and its function (e.g. Goldberg 2006). Yet, a great number, perhaps
 6 even a majority, of Construction Grammarians these days seem to take a usage-
 7 based approach to the study of patterns, which means that they consider a unit
 8 as a construction as soon as it has sufficient frequency (as evidenced by corpus
 9 data), regardless of whether or not that unit displays any sort of arbitrariness in
 10 the way its form links up with its function. This is also the approach taken here.
 11 We are less concerned with the potential unpredictability of a pattern's form or
 12 function than with its high frequency.

13 The structure of our paper is as follows. In Section 2, we will introduce the
 14 concept of n-grams. In Section 3 we will propose an application of n-grams
 15 to English language learning. Section 4 is devoted to some possible criticisms
 16 that could be levelled at this approach and to our rebuttal of them. Section 5
 17 presents some further features of an envisaged n-gram-based grammar, or 'n-
 18 grammar', of English, which is a project-in-progress. Our conclusions can be
 19 found in Section 6.

20

21

22 **2 What are n-grams, and what are they** 23 **typically used for?**

24

25

26 N-grams are sequences of n items, where n stands for any natural number (1, 2,
 27 3, 4, etc.) of linguistic units. For example, the word string *the fool on the hill*
 28 contains five 1-grams (usually called 'unigrams'), namely *the*, *fool*, *on*, *the* and
 29 *hill*, four 2-grams (or 'bigrams'), namely *the fool*, *fool on*, *on the* and *the hill*, three
 30 3-grams (or 'trigrams'), namely *the fool on*, *fool on the* and *on the hill*, two 4-
 31 grams, namely *the fool on the* and *fool on the hill*, and also one 5-gram, namely
 32 the string *the fool on the hill* itself. There are not just word-based n-grams but
 33 also character-based n-grams. Thus, the letter sequence *chat* consists of four
 34 unigrams (*c*, *h*, *a* and *t*), three bigrams (*ch*, *ha* and *at*), two trigrams (*cha* and
 35 *hat*) and one 4-gram (*chat*). The items in question that an n-gram has n adjacent
 36 instances of could be of any category. For instance, in Section 3, we will make
 37 use of n-grams whose items are word classes (determiner, noun, verb, etc.).

38

39

40

172 statistical language modelling. By 'language model', computational linguists

1 understand a set of probabilities (P's) which reflect, as accurately as possible,
 2 real language use. As Jurafsky (2012) puts it, “[i]t might have been better to
 3 call this ‘the grammar’. I mean, technically, what this is, is telling us some-
 4 thing about how [well] [...] words fit together, and we normally use the word
 5 ‘grammar’ for that, but it turns out that the word ‘language model’ [...] is
 6 standard”. Based on n-grams extracted from a large corpus, a language model
 7 may compute the likelihood of an entire string of n items (‘joint probability’)
 8 and/or the likelihood of a single upcoming item given $n-1$ previous items
 9 (‘conditional probability’). Estimates of these probabilities generated by an
 10 n-gram-based language model are used in a variety of practical applications.
 11 Table 1 gives some examples, drawn from Jurafsky (2012).

12
 13 **Table 1:** Some applications of an n-gram-based probabilistic language model (based on
 14 Jurafsky 2012)

15 Application	16 Task	17 Example
18 Machine translation	19 Distinguishing between ‘good’ 20 and ‘bad’ translations by their 21 probabilities	22 <i>High winds tonight</i> may be a 23 better translation than <i>large</i> 24 <i>winds tonight</i> , based on: 25 $P(\text{high winds tonight}) >$ 26 $P(\text{large winds tonight})$
27 Spell correction	28 Detecting likely mistakes based 29 on the probabilities of word 30 sequences	31 <i>The office is about fifteen minuets</i> 32 <i>from my house</i> likely contains a 33 misspelling from <i>minutes</i> , based 34 on: 35 $P(\text{about fifteen minutes from}) >$ 36 $P(\text{about fifteen minuets from})$
37 Speech recognition	38 Deciding between two sequences 39 that sound phonetically similar by 40 comparing their probabilities	41 <i>I saw a van</i> is likely to be a more 42 accurate transcription than <i>eyes</i> 43 <i>awe of an</i> , based on: 44 $P(\text{I saw a van}) \gg$ 45 $P(\text{eyes awe of an})$

31
 32 There are many other everyday applications. Word-based and character-based
 33 n-grams underlie features such as word suggestion and word completion avail-
 34 able on search engines and on our smartphones’ text messaging function.

35 While extracting n-grams from corpora is a common method of identifying
 36 recurrent formulae in discourse, other types of sequences are sometimes used
 37 apart from n-grams, such as so-called lexical bundles (Biber, Conrad, and Cortes
 38 2004), p-frames (Römer 2010) and skip-grams (Guthrie et al. 2006).

3 Using COCA n-grams for a new kind of grammar

3.1 The problem of ubiquitous constructions

Finding out what the most frequently used constructions are in a short text may sound like an easy enough task. In fact, it is not. To begin with, we would have to decide on an appropriate definition of ‘construction’. Secondly, suppose that we adopt a quite open definition of ‘construction’, as is common in Construction Grammar (e.g. Goldberg 2006), and count as construction every learned form-function pairing, ranging from individual words and morphemes to larger syntactic structures, it would then be hard not to overlook any of them, as any single sentence typically may contain one or several dozen constructions. This will become clear if we consider an example taken from Goldberg (2003):

(1) *What did Liza buy the child?*

This short sentence contains all of the following constructions:

- (2)
- a. the *buy, child, did, Liza, the* and *what* constructions (i.e. words)
 - b. the Ditransitive construction (i.e. double-object construction)
 - c. the Question construction (which is a fairly abstract construction, involving a certain intonation contour)
 - d. the Subject-Auxiliary Inversion construction (which is not only used in questions)
 - e. the VP construction
 - f. three cases of the NP construction (namely, *What, Liza* and *the child*)

For the time being, it is technically very hard, if not impossible, to detect and tally all these kinds of constructions automatically, which is what would be required if we wanted to count constructions in a whole corpus.

We propose to bypass the problem of scripting such a construction-detecting program by relying on readily available part-of-speech (POS) n-grams, which we will treat as constructions (or major parts thereof). This decision, of course, needs proper justification, which we will attempt to give in Section 4.1. At present, we are focusing on describing the methodology used.

3.2 The general idea

Via the website www.ngrams.info, one can download free lists of the most frequent 2-, 3-, 4- and 5-grams from COCA. Each list contains about 1,000,000 lexical n-grams. The lists are ordered from the most frequent to the least frequent n-grams. Table 2 gives some examples from the top of each list.

Table 2: N-grams from COCA, with some of the highest-frequency examples

N-grams	Examples
2-grams	<i>of the, is a, going to, I think, ...</i>
3-grams	<i>one of the, a lot of, the United States, as well as, ...</i>
4-grams	<i>I do-n't know, for the first time, on the other hand, ...</i>
5-grams	<i>I do-n't think so, the rest of the world, by the end of the, ...</i>

Observe, by the way, that the contracted negator (-n't) is treated as a separate word by the tagger.

Via the website mentioned above, it is also possible to download lists of n-grams where part-of-speech tags are presented together with the actual words making up each n-gram. What we claim here is that one can exploit this information to find common grammar structures in the corpus (and hence, to the extent that COCA is a representative corpus, in a major variety of the English language). For the purposes of illustration, Figure 1 shows the top section of the list of 4-grams containing part-of-speech information.

The left-most column gives us the number of occurrences ('tokens') of the lexical n-gram ('type') in question in COCA. Thus, *I don't know* is the most frequent 4-gram in COCA, occurring 54,632 times in the corpus, followed by *I don't think*, with 43,760 occurrences.

The four columns to the right contain the part-of-speech information, based on the CLAWS 7 tagset.⁴ Thus, the tag *ppis1* stands for 'singular personal pronoun, first person, subjective case' (i.e. the word *I*), *vd0* for 'do as a finite form (in declarative and interrogative clauses)', *xx* for 'not' or its contracted form, and *vi* for 'the base form of a lexical verb used as an infinitive'. As can be noticed, the first two n-grams have the same part-of-speech tags. They share this part-of-speech tagging with *I don't want*, a little further down the list (see the boxes with dotted lines). Similarly, the 4-grams *the end of the* and *the rest of the* (in 4th and 6th position) share their part-of-speech labelling (see the boxes with full lines). The idea now is to order all these POS 4-grams by their frequency,

⁴ <http://ucrel.lancs.ac.uk/claws7tags.html>, last accessed on 28 February 2015.

1	54632	I	do	n't	know	ppis1	vd0	xx	vvi
2									
3	43760	I	do	n't	think	ppis1	vd0	xx	vvi
4	33968	in	the	United	States	ii	at	np1	np1
5									
6	29848	the	end	of	the	at	nn1	io	at
7	27119	do	n't	want	to	vd0	xx	vvi	to
8									
9	21537	the	rest	of	the	at	nn1	io	at
10	19864	at	the	end	of	ii	at	nn1	io
11	19165	for	the	first	time	if	at	md	nnt1
12	18632	I	do	n't	want	ppis1	vd0	xx	vvi
13									
14	18115	at	the	same	time	ii	at	da	nnt1
15	16809	in	the	middle	of	ii	at	nn1	io
16	16681	one	of	the	most	mc1	io	at	rgt
17	16626	of	the	United	States	io	at	np1	np1
18	15857	is	one	of	the	vbz	mc1	io	at
19	14392	to	be	able	to	to	vbi	jk	to
20									
21									
22									
23									
24	...								
25									

Figure 1: Most frequent lexical 4-grams ('types') from COCA, together with their number of corpus occurrences ('tokens') and their part-of-speech tags

that is, by the number of different lexical 4-grams (lexical types) that instantiate them. What this reordering results in is shown in Figure 2.

The most frequent POS 4-gram in COCA is the one instantiated by *at the end of, in the middle of* and 6,984 other sequences of a preposition (other than *of*), the definite article, a singular common noun and the preposition *of*. Other than the list of lexical 4-grams (Figure 1), this list gives us direct information about what the most common syntactic structures are for 4-word sequences in COCA. We could use such a frequency list as the basis for an n-grammar of English. As an added bonus, we could combine this syntactic information with lexical information about the most common actual 4-grams (see also Section 3.4). Indeed, Construction Grammar assumes that both lexical chunks and the more general patterns they instantiate have their role to play in (first and second) language

1	6986	ii at nn1 io
2		
3	5382	nn1 io at nn1
4	4645	ii at jj nn1
5		
6	4235	nn1 ii at nn1
7		
8	4177	at jj nn1 io
9	3847	at nn1 io at
10		
11	3609	ii at1 jj nn1
12	3569	at1 jj nn1 io
13		
14	3313	to vvi at nn1
15	3249	at nn1 ii at
16		
17	3028	ii at nn1 nn1
18	2848	at nn1 io nn1
19		
20	2797	ii at nn1 cc
21	2684	at1 nn1 ii at
22		
23	2573	at1 jj nn1 ii
24	...	
25		

Figure 2: Most frequently instantiated POS 4-grams in COCA with number of lexical instantiations ('types') for each POS 4-gram

acquisition (cf. Ellis 1996, 2003, 2013; Tomasello 2003; see also Lewis 1993, the papers in Cowie 1998 and Wray 2002, inter alia, on the role of chunks in acquisition). Learners can be said to master a target language all the more accurately the more they manage to use lexical items in their preferred constructional environment (Wulff and Gries 2011).

Rather than using 2-, 3- or, as just demonstrated, 4-grams, we suggest using 5-grams as the basis of our n-grammar of English, which are the longest n-grams available from COCA's n-gram website. Traditional grammars tend to focus on shorter units, but if we want to target intermediate to advanced students, we believe that strings of 5 segments present an adequate size – neither too short, nor too long. While even longer n-grams could in principle have been used, if they had been available from COCA, we do not think such longer strings would

1 have provided many more relevant constructions, since longer strings are likely
 2 to be made up of shorter component structures, which we will show to be the
 3 case for 5-grams already. Moreover, in those cases in which a 5-gram does not
 4 coincide with a complete syntactic phrase, it can still be extended ‘by hand’
 5 with a phrasal category, something which will also be illustrated below. In short,
 6 our choice of using complete or extended 5-grams is thus motivated by the aim
 7 to use units that are as long and complete as possible, but which at the same
 8 time still allow manipulation and combination to form even larger structures in
 9 the language.

12 3.3 The method in detail

13 We restrict our selection to the 100 most frequent POS 5-grams based on the
 14 COCA list of lexical 5-grams containing part-of-speech information. In an
 15 n-gram-based grammar of English, each such pattern could and should be
 16 presented together with some of its frequent lexical instantiations, so as to
 17 show how the skeletal structures can be fleshed out in actual language use.
 18 Why 100 patterns? This is a somewhat arbitrary choice, motivated less by lin-
 19 guistic factors than by reasons related to learner motivation: learners might
 20 consider 100 patterns an achievable target. Needless to say, one could also
 21 select 200 patterns, 500 patterns, etc., or alternatively 365 patterns, one for each
 22 day of the year.

23 The list of lexical 5-grams with part-of-speech tags that can be downloaded
 24 from Mark Davies’s website mentioned above (www.ngrams.info) contains
 25 exactly 1,293,537 types of lexical strings. The list is cut off at 5-grams with a
 26 minimum frequency of 5 occurrences in the corpus (presumably because the
 27 list was meant to contain ca. one million types). Remember from Table 2 that
 28 this list contains such sequences as *I don’t want to* or *the rest of the world*. We
 29 grouped these lexical strings according to the syntactic patterns they instantiate
 30 (i.e., their part-of-speech tag sequence). We thus obtained a total of 325,552 POS
 31 5-grams. The number of lexical strings (‘types’) per POS 5-gram varies from 7,272
 32 for *at nn1 io at nn1*, a structure shared by *the rest of the world*, *the side of the*
 33 *road* and thousands more (where *at* stands for *the*, *nn1* for a singular common
 34 noun and *io* for *of*), to just 1, for instance in the case of *appge cc appge jjt nn1*
 35 for *his or her best interest* (in which *appge* stands for a possessive determiner, *cc*
 36 for a coordinating conjunction and *jjt* for a superlative adjective).

37 We then took the 100 most frequent POS 5-grams as main content for the
 38 n-grammar. By ‘most frequent’ POS 5-grams, we mean those that represent the
 39 highest number of types, that is, the highest number of different lexical 5-grams
 40

1 that have the structure specified by them. These patterns (for which we also
 2 made the part-of-speech labels more transparent) now range in frequency from
 3 still 7,272 types for the sequence [*the* X_{noun} *of the* Y_{noun}] to 499 types for the
 4 sequence made up of a complex 3-word preposition followed by *the* and an
 5 adjective (e.g. *in front of the whole*).

6 As is clear from this last example, a 5-gram does not necessarily form a
 7 complete constituent, since n-grams are ‘blind’ to constituent structure. Some-
 8 times, an n-gram does not contain enough (or one might say, it may contain
 9 too much) to make up what we would intuitively consider an ordinary linguistic
 10 sequence. We therefore added an element to the right in those cases where the
 11 right-most boundary of a 5-gram does not coincide with a closing bracket, so to
 12 speak. Thus, in the case of the pattern [*at/in/to/... the* X_{noun} *of the*], we just add
 13 an element to the right of the determiner. This element could be a noun, but it
 14 could also be an adjective which precedes a noun, among other possibilities.
 15 Technically, the grammatical category covering all of these is what is called
 16 a ‘nominal’ (nom) in Huddleston, Pullum et al.’s (2002) grammar, or an ‘N-bar’
 17 (N[̄]) in X-bar theory (Chomsky 1970; Jackendoff 1977) – that is, a noun phrase
 18 minus the determiner. We always completed with a category label that stands
 19 for the widest range of possible continuations. Additions are between parentheses
 20 in what follows. Table 3 shows the result of the procedure for the ten most
 21 common POS 5-grams, along with an example of each.

22 **Table 3:** Ten most frequent syntactic (completed) 5-grams in COCA
 23

Syntactic pattern	Example
<i>the</i> X_{noun} <i>of the</i> Y_{noun}	<i>the rest of the world</i>
<i>at/in/to/... the</i> X_{noun} <i>of the</i> (Y_{nom})	<i>at the end of the (day)</i>
X_{noun} <i>at/in/to/... the</i> Y_{noun} <i>of</i> (Z_{NP})	<i>increase in the number of (students)</i>
<i>the adj</i> X_{noun} <i>of the</i> (Y_{nom})	<i>the other side of the (room)</i>
<i>to verb the</i> X_{noun} <i>of</i> (Y_{NP})	<i>to improve the quality of (life)</i>
<i>at/in/to/... the adj</i> X_{noun} <i>of</i> (Y_{NP})	<i>on the other side of (the room)</i>
<i>a(n) adj</i> X_{noun} <i>at/in/to/... the</i> (Y_{nom})	<i>a far cry from the (original proposal)</i>
<i>at/in/to/... the</i> X_{noun} <i>of one's</i> (Y_{nom})	<i>at the top of his (lungs)</i>
<i>the</i> X_{noun} <i>of the adj</i> (Y_{nom})	<i>the end of the Cold (War)</i>
<i>at/in/to/... the</i> X_{noun} <i>at/in/to/... the</i> (Y_{nom})	<i>on the way to the (hospital)</i>

36 We never added any symbol to the left of the 5-gram, as the left-most symbol
 37 always constitutes the first element of a constituent. However, in some cases, a
 38 lexical instantiation of a POS n-gram might benefit from one or more added ele-
 39 ments at the left. For instance *the benefit of the doubt* is a syntactically complete
 40

1 unit – it is an NP – but it typically occurs in lexically larger environments,
2 involving verbs such as *give*, *get* or *deserve*.

3.4 The medium-level and hybrid nature of the POS n-grams

6 As the reader will have noticed, the syntactic information associated with the
7 n-grams from COCA consists of quite specific part-of-speech tags. While there
8 are only eight or nine word classes traditionally recognized in English, the
9 CLAWS 7 tagset, which was used for tagging COCA, distinguishes between 137
10 different categories. For instance, the preposition *of* is treated not as any pre-
11 position but as the word *of*, all the forms of the verb *be* are treated differently
12 from each other, singular nouns are treated as different from plural nouns, com-
13 mon nouns as different from proper nouns, and so on. As a result, the POS
14 n-grams are not maximally general, as would have been the case if they were
15 of the type ‘Det N Prep Det N’. Nor of course are these POS n-grams maximally
16 specific, which is the case only for purely lexical n-grams such as *the name of*
17 *the motel*, where each item is an actual word. While it would be possible for us
18 to come up with more general patterns based on the specific part-of-speech tags
19 and to calculate their frequencies, our POS n-grams as they are may in fact come
20 close to having the ideal grain size of a construction: neither too schematic nor
21 too concrete. We do not want to claim here that there are no such things as
22 very general constructions or that some specific lexical strings cannot have the
23 status of stored language units; indeed, a standard assumption of Construction
24 Grammar is that generalizations over exemplars and the (sufficiently frequent)
25 exemplars themselves are stored in the speaker’s mind (cf. Section 4.1). Yet,
26 while it is obvious that specific items *have to* be stored if they are formally or
27 semantically unpredictable – storage is required for words and idioms – there
28 is no equally compelling reason why we would need to store the most schematic
29 rules of language. As Croft (1998: 168) formulates it, “[s]peakers do not neces-
30 sarily make the relevant generalizations, even if clever linguists can”.⁵

33 ⁵ In actual fact, Croft’s (1998) quoted sentence is lifted from an article that deals more with
34 semantics (polysemy and homonymy) than with the level of generality at which speakers store
35 constructions. However, these issues are not unrelated and lie at the heart of the difference
36 between a Goldbergian (1995) approach to argument structure constructions (i.e. one in which
37 they are treated as highly schematic form-meaning patterns) and a Boas-style (2003) or Iwata-
38 style (2008) approach to them (i.e. one in which so-called ‘mini-constructions’ or specific lexical
39 constructions are associated with individual verbs or even individual verb senses). See also
40 Levshina and Heylen (2014) for related findings about the optimality of medium-level granularity
in the context of semantic classes of predicates governing the choice between competing
constructions.

1 Apart from, or as a corollary of, being somewhat below the maximum level
 2 of generality, our POS n-grams are also somewhat hybrid in nature, that is, they
 3 are *partially* rather ‘syntactic’ and *partially* rather ‘lexical’. The mixing of levels
 4 results purely from the rich, fine-grained tagset that is used for the COCA corpus,
 5 but there are computationally more sophisticated methods for automatically
 6 generating linguistically ‘interesting’ n-grams which combine lexical items and
 7 formal categories: see Wible and Tsao (2010), Lyngfelt et al. (2012) and Forsberg
 8 et al. (2014). Again, by mixing more general and more specific items in a single
 9 template, we may approximate the ideal of constructions viewed as language
 10 units that actually operate in the mind of speakers. For instance, in our top
 11 hundred POS n-grams, we find the following hybrid structures:⁶

- 12
 13 (3) a. *the* X_{noun} *of the* Y_{noun} (pattern No. 1), e.g. *the rest of the world*
 14 b. *the* X_{noun} *of a(n)* Y_{noun} (pattern No. 20), e.g. *the son of a bitch*
 15 c. *a(n)* X_{noun} *of the* Y_{noun} (pattern No. 30), e.g. *a thing of the past*
 16
 17

18 The trained linguist may need some convincing to see that these are distinct
 19 patterns and not *just* different realizations of a single more general pattern. But
 20 notice, first of all, the difference in frequency. The POS n-gram in (3a), as we
 21 noted above, covers 7,272 lexical types (with at least five tokens, i.e. corpus
 22 occurrences, each), while the ones in (3b) and (3c) only cover 1,448 and 968
 23 lexical types, respectively. The pattern [*a(n)* X_{noun} *of a(n)* Y_{noun}], with two in-
 24 definite articles, does not even rank among the hundred most frequent POS
 25 5-grams. Secondly, while each pattern provides open slots for nouns, they do
 26 not allow the same nouns in these slots. For instance, we would not find *?the*
 27 *son of the bitch*, *??the thing of the past* or *??a rest of the world*. This suggests
 28 that each pattern has its own particular properties, causing it to attract certain
 29 nouns and to repel certain others. We will come back to this in Section 5.3.

30 Because of a pattern’s close association with *some* lexical items and not
 31 with others, we feel that it is worthwhile to provide this information to learners.
 32 This is fully in line with a constructionist and usage-based approach to lan-
 33 guage learning, which stresses the importance of exemplars in acquisition
 34 (Abott-Smith and Tomasello 2006; Ellis 2006, 2013). So, ideally, an n-grammar
 35 should present not just semi-schematic and hybrid patterns but also some of
 36

37
 38 ⁶ The first segment is not always the determiner *the* but could also be the quantifier *no*, as in
 39 *no mention of the fact (that...)*. Because this determiner is used far less frequently than *the*, we
 40 use the latter as a transparent substitute for the tag ‘at’. A rare example in which the quantifier
no is used before the second noun is *the point of no return*.

1 the frequent lexical instantiations that they generalize over. In the case of the
 2 pattern [*the* X_{noun} *of the* Y_{noun}], this would mean that the learner also gets
 3 to see some fully lexical sequences, possibly even with corpus frequencies
 4 (number of tokens) added to them, as shown in Table 4.

5
 6 **Table 4:** Most frequent lexical realizations of the pattern
 7 [*the* X_{noun} *of the* Y_{noun}] in COCA

8	<i>the rest of the world</i> (3,618)	<i>the benefit of the doubt</i> (547)
9	<i>the side of the road</i> (1,217)	<i>the edge of the bed</i> (530)
10	<i>the rest of the country</i> (1,174)	<i>the center of the room</i> (526)
11	<i>the fact of the matter</i> (825)	<i>the State of the Union</i> (495)
12	<i>the end of the world</i> (764)	<i>the back of the room</i> (463)
13	<i>The fact of the matter</i> (717)	<i>the back of the head</i> (450)
14	<i>the end of the war</i> (670)	<i>the middle of the room</i> (448)
15	<i>the rest of the way</i> (597)	...

16 Note, incidentally, that *the fact of the matter* appears twice. This is due to the fact
 17 that the downloadable lexical n-grams with POS information are case-sensitive,
 18 which means that a word with a capital letter and the same word without
 19 a capital letter are treated as belonging to different n-grams. This may seem
 20 like a nuisance, but it actually provides useful information about where that
 21 n-gram is found in the sentence (sentence-initially or not).

24 4 Possible points of criticism and their rebuttal

26 We are aware of some immediate objections that one might raise against the
 27 approach we take to selecting patterns to be included in a new, radically
 28 usage-based grammar. We can think of at least the following four points of
 29 criticism:

- 30
 31 (i) Not all of these POS 5-grams are constructions.
 32 (ii) There is a lot of (and perhaps too much) overlap between them.
 33 (iii) The top hundred POS 5-grams are but the tip of the iceberg.
 34 (iv) By restricting ourselves to 5-grams, we may miss out on interesting 2-, 3- and
 35 4-grams.

36
 37 In the following subsections, we will defend our approach against this possible
 38 criticism.

40

4.1 Not all of them constructions?

One might wonder what is so special about, for example, the pattern [*the* X_{noun} *of the* Y_{noun}], which appears to be formed on the basis of some general phrase-structure rules, namely the ones listed in (4), combined with the knowledge that *the* is a determiner and *of* a preposition:

- (4) a. NP → det nom
 b. nom → noun (PP)
 c. PP → prep NP

It is true that on Goldberg's (1995: 4) original definition of the term, this first pattern and many (if not most) of the other patterns in our top hundred POS 5-grams would not qualify as constructions. This definition stated that a form-function pairing is a construction only if there is something about its form or function that is not strictly predictable from what is already available in the grammar. However, in Goldberg's (2006: 5) later work, this requirement is loosened: "Any linguistic pattern is recognized as a construction as long as some aspect of its form or function is not strictly predictable from its component parts or from other constructions recognized to exist. *In addition, patterns are stored as constructions even if they are fully predictable as long as they occur with sufficient frequency*" [emphasis ours]. The view that constructions are psychologically entrenched form-function pairs is also expressed by Croft and Cruse (2004: 288), Langacker (2005: 140) and Bybee (2006: 715).

Given their high frequency, which is the basis of their selection in the first place, there is no doubt that all of the hundred POS 5-grams meet this definition of 'construction': even if some of them are compositional, their sheer frequency makes it unlikely that they are formed anew each time they are used. It is much more plausible that these patterns are directly retrieved from what construction grammarians, using a term coined by Jurafsky (1991), call the 'construct-i-con'.

4.2 Too much overlap between them?

It will have been noted that many of the patterns shown in Table 3 look like variations on a theme. That is, they do not all represent fully distinct constructions. For example, the two most frequent POS 5-grams ([*the* X_{noun} *of the* Y_{noun}] and [*at/in/to/... the* X_{noun} *of the* (Y_{nom})]) share four fifths of their component elements (namely '*the* X_{noun} *of the*').

1 As we see it, however, this overlap is not a problem. On the contrary, it
 2 allows us to integrate in this new type of learner's and student's grammar an
 3 important feature of language: the possibility of reusing parts of structures over
 4 and over again in slightly different environments or with slight modifications. In
 5 other words, one of the main properties of language is that it involves structures
 6 that are partially reusable and adaptable. This feature can be visualized in chart
 7 form, which for each pattern shows its relatedness to some of the other patterns
 8 in the grammar. Figure 3 is an example of what such a chart could look like for
 9 the first pattern in the n-grammar.

10

11 No. 1		<i>the</i>		X_{noun}	<i>of</i>	<i>the</i>		Y_{noun}
12 No. 2	<i>at/in/to/...</i>	<i>the</i>		X_{noun}	<i>of</i>	<i>the</i>		Y_{nom}
13 No. 4		<i>the</i>	<i>Adj</i>	X_{noun}	<i>of</i>	<i>the</i>		Y_{nom}
14 No. 9		<i>the</i>		X_{noun}	<i>of</i>	<i>the</i>	<i>adj</i>	Y_{nom}
15 No. 13		<i>the</i>		X_{noun}	<i>of</i>	<i>one's</i>		Y_{noun}
16 No. 18		<i>the</i>		X_{noun}	<i>at/in/to/...</i>	<i>the</i>		Y_{noun}
17 No. 20		<i>the</i>		X_{noun}	<i>of</i>	<i>a</i>		Y_{noun}

18

18 **Figure 3:** 'Chop and change' chart for the first pattern in the n-grammar (top row)

19

20

21 The chart shows how a pattern can be chopped up to allow for the insertion of
 22 elements in the right position (e.g. adjectives before nouns) and be changed by
 23 replacing elements in a structural position by alternatives in that position (e.g.
 24 an indefinite article or a possessive determiner instead of a definite article). As
 25 such, this 'chop and change chart' (a term whose rights of use in grammar
 26 instruction we hereby reserve) directly represents the syntagmatic and para-
 27 digmatic relations between structural elements in the grammar. The overlap
 28 between patterns reflects the fact that grammar is a combinatorial system, which
 29 operates on classes of discrete string segments such as adjectives or nouns.
 30 This is not just how linguists see it, but how the human brain treats grammar
 31 (Pulvermüller and Knoblauch 2009; Pulvermüller, Cappelle, and Shtyrov 2013).

32

33 The n-grammar proposed here is thus not a maximally parsimonious system
 34 to generate word sequences. Instead, it represents structural information in a
 35 way that is full of redundancy. This redundancy may be helpful for learners to
 36 master the structures: though they could have been captured more economi-
 37 cally, this would have been at the expense of a range of learning opportunities
 38 spread out through time, which is needed for consolidated acquisition. This is as
 39 true for humans as it is for the simplest of organisms. To cite the Nobel prize
 40 winning neuroscientist Eric R. Kandel: "Conversion of short-term to long-term

40

1 memory storage requires spaced repetition – practice makes perfect, even in
2 snails” (Kandel 2001: 1031).

3 4 5 **4.3 Just the tip of the iceberg?**

6 Remember that for the n-grammar proposed here, we retained a mere 100 POS
7 5-grams out of 325,552 such patterns in COCA (themselves generalizing over
8 more than a million lexical 5-grams with at least 5 occurrences each in the
9 corpus). Put differently, our top hundred most frequently occurring syntactic
10 templates only represent 0.03% of all possible POS 5-grams based on lexical
11 5-grams with at least 5 tokens in COCA. Put differently still, our selection does
12 not seem to count for much. Using familiar imagery, if the part of an iceberg
13 that appears above water is only one tenth of its total volume, then our selection
14 does not represent the proverbial tip of the iceberg, and not even the tip of the
15 tip of the tip of the iceberg.

16
17 If we disregard POS 5-grams that are instantiated by fewer than 5 different
18 lexical 5-grams (‘types’), there are no longer 325,552 POS 5-grams in COCA, but
19 36,617 of them. Our selection then represents 0.27% of these. This is admittedly
20 still a very small portion, which is barely visible in the left-most stacked bar in
21 Figure 4. However, if we now look at the types represented by these 36,617 POS
22 5-grams, we find that there is a total number of 823,683 different lexical 5-grams
23 in COCA. Our top hundred POS 5-grams represent a non-negligible portion of
24 these: 105,184 types, which is 12.77% (cf. middle stacked bar in Figure 4). In
25 terms of tokens (individual occurrences), there are 11,248,178 sequences of 5
26 words corresponding to the 36,617 POS 5-grams. Our top hundred POS 5-grams
27 cover 1,615,199 tokens. This high number represents 14.36% of the total number
28 of tokens for all 36,617 POS 5-grams (cf. the stacked bar on the right in Figure 4).
29 What we find here is something akin to what Zipf (1935) noted for lexical items,
30 namely that the most frequent items (types) cover a large part of the occurrences
31 (tokens) in usage. In the Brown corpus, for instance, half of the word volume is
32 accounted for by only 135 vocabulary items (Fagan and Gençay 2010).

33 In sum, even our very small set of POS 5-grams (just one hundred out of
34 more than thirty thousand in the corpus) appears to have quite large coverage
35 in terms of number of lexical strings (both the types and their tokens) that
36 correspond to these syntactic templates. It exceeds the percentage of an iceberg
37 that extends above the water surface.

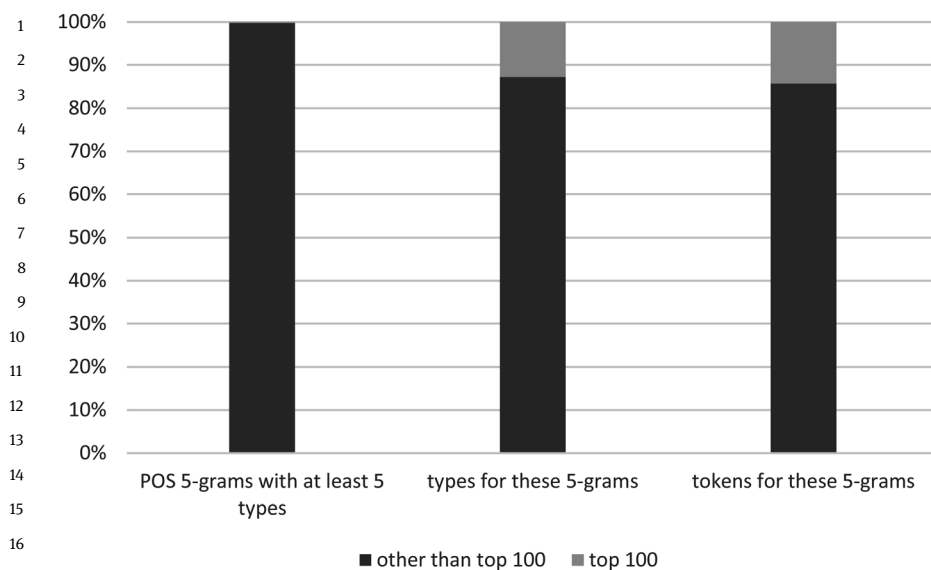


Figure 4: Share of the selected highest-frequency POS 5-grams in COCA and its coverage in terms of types and tokens

4.4 Neglecting important 2-, 3- and 4-grams?

It may seem an odd choice to take 5-grams as our basis for a grammar of English. Even quite apart from the fact that some constructions are longer than five words or may involve discontinuities and hence cannot be captured as 5-grams, there is the risk of overlooking interesting grammar patterns that are shorter than 5 words. For instance, among the top ten most frequently used 2-grams, we found the pattern $[X_{\text{noun}} Y_{\text{noun plur}}]$, that is, a noun-noun compound in the plural, such as *family members*, *interest rates*, *phone calls*, *college students*, *side effects*, and several thousands more. Unfortunately, the two-word POS sequence $[X_{\text{noun}} Y_{\text{noun plur}}]$ is not a part of any of our top hundred POS 5-grams. Does this example not suggest that we may fail to integrate some vital grammar patterns by focusing on 5-grams only?

While any common lower- n -gram pattern that is not included in any of our 5-gram patterns is surely a missed opportunity to capture all that is essential in grammar, the actual situation does not give reason for too much concern. Table 5 shows that the five highest-frequency POS 2-grams taken together are included 102 times in our selection of POS 5-grams, and that the five highest-frequency POS 3-grams and POS 4-grams are still also included 64 times and 35 times, respectively. (Obviously, shorter sequences have a higher likelihood of being included than longer ones.)

Table 5: Inclusion of high-frequency POS 2-, 3- and 4-grams from COCA in the hundred most frequent POS 5-grams from COCA

N-grams with $n < 5$	Top 5 most frequent POS n-grams	Number of inclusions in top 100 POS 5-grams
2-grams	1. adj X_{noun}	19
	2. $X_{\text{noun}} Y_{\text{noun}}$	3
	3. adj $X_{\text{noun plur}}$	2
	4. $X_{\text{noun}} \{at/in/to/.. \}$	20
	5. <i>the</i> X_{noun}	58
3-grams	1. $\{at/in/to/.. \}$ <i>the</i> X_{noun}	27
	2. <i>a(n)</i> adj X_{noun}	9
	3. <i>the</i> adj X_{noun}	7
	4. $X_{\text{noun}} \{at/in/to/.. \}$ <i>the</i> (X_{nom})	14
	5. $\{at/in/to/.. \}$ <i>an</i> X_{noun}	7
4-grams	1. $\{at/in/to/.. \}$ <i>the</i> X_{noun} <i>of</i> (Y_{NP})	17
	2. X_{noun} <i>of the</i> Y_{noun}	6
	3. $\{at/in/to/.. \}$ <i>the</i> adj X_{noun}	3
	4. $X_{\text{noun}} \{at/in/to/.. \}$ <i>the</i> Y_{noun}	6
	5. <i>the</i> adj X_{noun} <i>of</i> (Y_{noun})	3

Clearly, it is not the case that by looking at 5-grams only, we ignore <5-grams. As explained in Section 2, a 5-gram by its very nature simultaneously harbors two 4-grams, three 3-grams and four 2-grams. As a consequence, if we study a hundred 5-grams, we actually get to see a thousand n-grams, not even counting the individual words. This is not to say that we get nine hundred *different* 2-, 3- and 4-grams for free with our 5-grams, as many of these lower-n-grams will be included several times. This is not a problem, in light of our discussion of redundancy and repetition in Section 4.2.⁷

⁷ To look at this from a somewhat different perspective, we might say that the high frequency of certain 5-grams accounts to some extent for the frequency of some of its included <5-grams. This is what O'Donnell (2011) points out to be the case for lexical n-grams: for instance, at the top of the list of lexical 5-grams in COCA we find *I don't want to*. This occurs 12,659 times in the corpus and thereby contributes in no small way to the high frequency of its multiple component parts (*I, do, n't, want, to, don't, want to, I don't want, etc.*). Therefore, by studying high-frequency 5-grams, the learner is given a glimpse into some of the reasons why smaller combinations are so frequent. This is true, we feel, for both n-gram templates of the sort discussed in our text and lexically 'filled-in' n-grams, of the sort O'Donnell (2011) focuses on.

1 By selecting 5-grams, we automatically retrieve more complex structures.
 2 This is why a ‘5-grammar’ of English may be more ideally suited for intermediate
 3 to advanced learners of English than for absolute beginners. For lower-level
 4 learners, n-grams other than 5-grams (namely, 2-grams, 3-grams and 4-grams)
 5 might be a better way to start. In other words, we do not want to claim that
 6 using 5-grams is the only valid way of constructing an n-grammar of English.

9 5 Further features of the n-grammar

11 5.1 Adding a visual measure of productivity

13 Pedagogical grammars do not generally contain any statistical information
 14 about frequency, unlike modern dictionaries, many of which provide an indica-
 15 tion of how common a word is, or in which genre or register it is typically used.
 16 Our proposed n-grammar can easily include such information. In Sections 3.2
 17 and 3.4 we already suggested that individual lexical n-grams associated with
 18 the more schematic POS n-grams may be shown with their actual corpus fre-
 19 quencies, thus giving the learner some idea of their usefulness as chunks in the
 20 target language. We believe that if learners see, for example, that *The fact of the*
 21 *matter* at the beginning of a sentence is used more than 700 times in a corpus of
 22 native-speaker English, this kind of knowledge may cause them to take note
 23 of this expression more consciously and stimulate them to use it themselves
 24 (cf. e.g. Schmidt (1990) and Robinson (2006) on the ‘noticing’ hypothesis and
 25 the role of conscious attention in second language acquisition). But the patterns
 26 themselves could also be provided with frequency information. Thus, the most
 27 frequent POS 5-gram, [*the X_{noun} of the Y_{noun}*], might be stated explicitly to have
 28 7,272 types (with at least 5 tokens each). In addition, we might mention that
 29 these types together represent 126,077 tokens in the corpus.

30 Such figures may not mean much by themselves to the learner. Though such
 31 high numbers might of course be impressive and therefore encourage the
 32 learner to devote due attention to the pattern, they will vary from corpus to
 33 corpus. A more general indication of frequency, similar to what can be found in
 34 certain dictionaries (e.g. high-frequency, medium-frequency, low-frequency)
 35 could be sufficient, if it were not for the fact that the top hundred most frequent
 36 POS 5-grams are naturally all at the high end of the frequency scale anyhow. It
 37 would probably be more beneficial to the learner to have a direct visual indica-
 38 tion of a pattern’s *usefulness*. If by ‘useful’ we mean how many different lexical
 39 realizations the pattern allows the learner to form, we should include a measure
 40 of the pattern’s productivity. Productivity can be defined in terms of the ratio of

1 types per tokens: the more types per number of tokens, the more productive a
 2 pattern is. This is clear if we consider the other extreme case: if all the corpus
 3 occurrences of a pattern were instances of only one lexical string, that ‘pattern’
 4 would have no productivity at all. Alternatively, productivity can be expressed
 5 in terms of the ratio of unique corpus occurrences (‘hapax legomena’, i.e.
 6 types with only one token) per tokens (cf. Baayen 1989): the more such single-
 7 occurrence types, the higher the probability that also ‘outside’ the corpus the
 8 pattern will be used to form novel creations and so the more productive the
 9 pattern. We propose here to combine the two measures (type-to-token ratio and
 10 hapax-to-token ratio) in a single graph.

11 There is one slight problem to overcome. Remember that all the lexical
 12 n-grams used for our n-grammar have at least five corpus occurrences, so that,
 13 strictly speaking, there are no hapax legomena among them. Therefore, we need
 14 to rely on a related statistic, which we could call ‘pentakis legomena’, that is,
 15 sequences that occur only five times in the corpus. The ratio of these, too, just
 16 like the ratio of hapax legomena, can give us an idea of how readily novel com-
 17 binations are formed based on a given pattern. The cut-off of five occurrences
 18 per type (cf. Section 3.3) also results in a somewhat skewed type/token ratio:
 19 above this cut-off point, there is a smaller type/token ratio (as here we find types
 20 with comparatively many tokens) than below that cut-off point (where we find
 21 types with relatively few tokens). Our solution to compensate for this skewing
 22 is to multiply the type/token ratio by 5. This makes mathematical sense: suppose
 23 all the lexical types had just five occurrences, then the unadjusted type/token
 24 ratio would be 0.2, and by multiplying this by 5, we would obtain the maximal
 25 productivity score of 1, which would be just what we would like to find in that
 26 situation. Likewise, the pentakis/token ratio is also skewed compared to the
 27 more commonly used hapax/token ratio, since for any pattern, if there are many
 28 ‘pentakis’, one could expect there to be even more hapaxes. So, there is natu-
 29 rally a smaller ‘pentakis’/token ratio given a cut-off restriction of 5 occurrences
 30 than there would be a ‘hapax/token’ ratio if the cut-off restriction was removed
 31 (and this is so even if the total number of tokens would of course also increase if
 32 we removed the cut-off restriction). The solution, here too, exists in multiplying
 33 the pentakis/token ratio by 5. The formula for the combined and adjusted type/
 34 token and pentakis/token measure of productivity of a pattern is given below,
 35 whereby n stands for the total number of lexical types instantiating the pattern
 36 with at least five occurrences, N the total number of tokens for all these types
 37 and p the number of lexical types with just five occurrences:

38
 39
 40

$$Productivity = \frac{\left(\frac{n}{N}\right) \cdot 5 + \left(\frac{p}{N}\right) \cdot 5}{2}$$

1 What this says is that the productivity of a pattern can be calculated by taking
 2 the average of its type/token ratio multiplied by five and its pentakis/token ratio
 3 multiplied by five. Thus, for the pattern [*the* X_{noun} *of the* Y_{noun}], the number of
 4 lexical types n is 7,272 types, the total number of corpus occurrences N is
 5 122,685 and the number of pentakis legomena p is 1,670. If we feed these
 6 numbers in the formula above, we get the following result:

$$7 \quad \text{Productivity} = \frac{\left(\frac{7,272}{122,685}\right) \cdot 5 + \left(\frac{1,670}{122,685}\right) \cdot 5}{2} = 0.186$$

11 This result can be represented in graph form on a scale from 0 to 1. A theoretical
 12 zero value of productivity would be obtained for a pattern where all types are
 13 prefabricated chunks. The value 1, for full productivity, would be the score for
 14 a pattern where all types are novel creations (or at least, where they are all
 15 pentakis). For ease of visual interpretation and comparison with other patterns,
 16 we use a logarithmic scale of 10, with the minimum value approximating zero
 17 and the maximum value 1.

18 Figure 5 charts the productivity of the first and second most frequent POS
 19 5-gram in COCA. The lower productivity of [*at/in/to/... the* X_{noun} *of the* (Y_{nom})]
 20 is explained by the fact that the part (Y_{nom}) is not actually included in the
 21 5-gram and so plays no role in the type and token data used for the calculation.
 22 This ‘extended’ 5-gram thus contains only two open slots, for a preposition and
 23 a noun, only the latter of which is an open word class. In the first pattern, there
 24 are two slots for an open word class, so the productivity of this pattern is
 25 obviously much higher. One might wonder whether it makes sense to add items
 26 that are not taken into account when counting the number of occurrences. The
 27 reason why we did this is that we want to show learners how a POS 5-gram can
 28 be used grammatically. If a POS 5-gram ends in a determiner or an adjective, we
 29 find it useful to state what the next element will be (a nominal). It should be
 30 obvious, though, that this element cannot be taken into consideration when we
 31 want to compare POS-grams for frequency, as there is no easy way to list up all
 32 the possible instantiations of this element, which could be a bare noun, a noun
 33 preceded by one or more adjectives, a noun followed by a prepositional phrase
 34 of any length, and so on.

35 It is important to make the learner see that a low-productivity pattern does
 36 not necessarily equal an uninteresting one. The lower the productivity, the
 37 greater the role of strongly entrenched sequences, which are responsible for the
 38 high token frequency. Thus, while this second pattern is clearly less productive,
 39 its most common type ([*at the end of the* Y_{nom}]) has 10,663 occurrences in COCA,
 40 against only 3,618 occurrences for the most frequently used lexical sequence

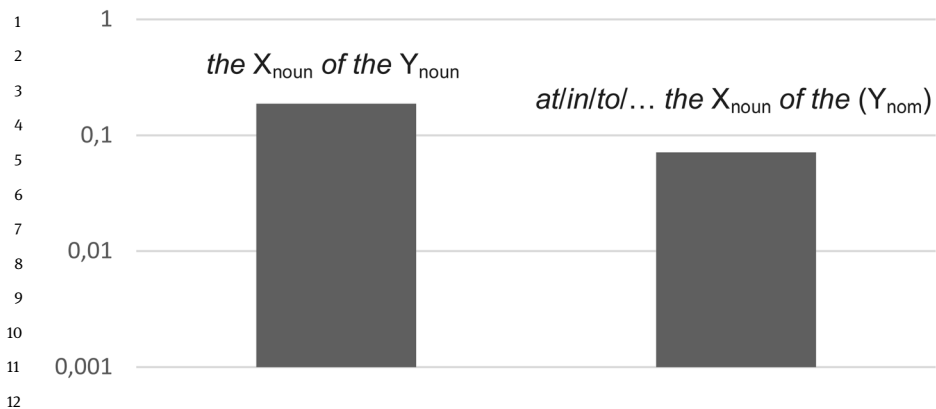


Figure 5: Productivity of two grammar patterns, visualised on a logarithmic scale

instantiating the first pattern (*the rest of the world*). To avoid the automatic association of higher with ‘better’, a suitable alternative visual representation might be one that plots the productivity score on a horizontally-oriented scale, where patterns towards the left margin are more ‘chunkified’ or ‘lexical’ and patterns towards the right margin are more ‘gridlike’ or ‘syntactic’.

5.2 Adding a visual measure of dispersion

5-grams may not appear as frequently in some genres or registers as they do in others. To indicate how evenly or how skewed a language item appears in different sections of a corpus, we can (or even should) use measures of dispersion (cf. Gries 2008a). Figure 6 illustrates a visually attractive way of showing which of the large components of COCA make use of the pattern [*the X_{noun} of the Y_{noun}*] the most and the least. The data were obtained by entering the query ‘the [*nn1*] of the [*nn1*]’ in the COCA search interface and looking up how many hits we retrieve in each of the main components of the corpus (spoken, fiction, magazine, newspaper and academic). This is information which the COCA search interface provides at a click of the mouse. Note that this search retrieves results for types whose token frequency is also lower than five. As could be expected of a rather complex NP, we find this grammatical structure used least frequently in spoken English and most frequently in academic writing. The spread through the corpus is clearly uneven.

Besides making immediately clear the uneven frequency of the pattern across broad corpus components, the graph can be adapted to indicate in which of these components, if any, the pattern in question occurs *much* more/less fre-

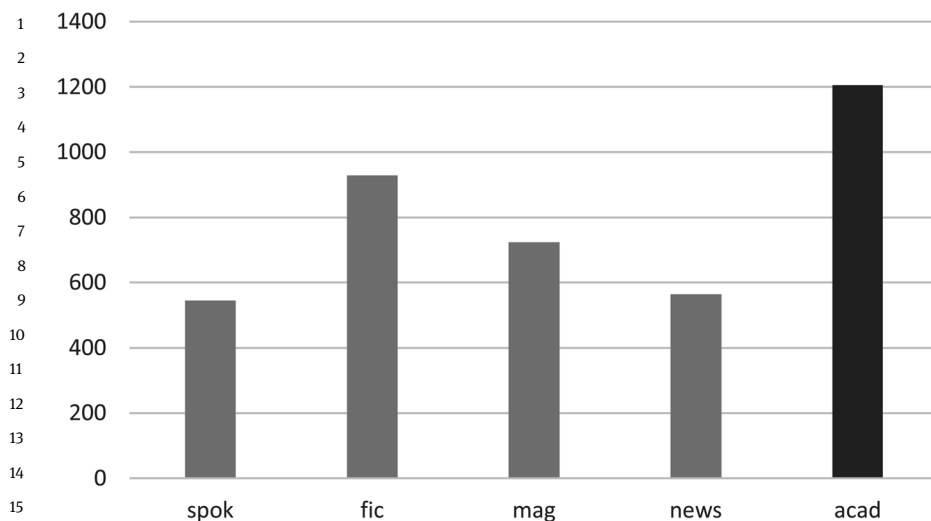


Figure 6: Dispersion of the X_{noun} of the Y_{noun} across COCA components. Bars indicate number of tokens per million words.

quently than expected (given an average). Specifically, we can use darker grey and lighter grey for markedly higher or lower frequency, respectively. This allows the learner to quickly identify the genre(s) in which the pattern is more conspicuously present or absent than what could have been expected under the assumption that it had an equal chance of occurrence in each component. We here define a markedly higher or lower frequency as a difference of at least 50% (for a positive difference, i.e. a surplus) or at least 33.33% (for a negative difference, i.e. a shortage) compared to the expected frequency. In this case, the expected frequency is 790 hits per million words (a figure obtained by dividing the total number of occurrences by the total corpus size multiplied by one million). Only in academic written discourse does the pattern display a marked difference (namely, an overuse of 53%) between what is observed and what is expected. In spoken discourse, the pattern is not sufficiently underrepresented – there is an underuse of (only) 31% – for its relative infrequency in this component to be considered of significance to the learner. This is why only the bar corresponding to academic writing has been given a different grey shade in Figure 6.

This uneven dispersion suggests that if we had used another corpus (containing for instance no academic writing at all), the frequency ranking of our POS n-grams might have been very different. The same may be true if we had used a corpus representing another variety of English. Obviously, corpus results

1 depend on (and vary with) the corpus used. This is also valid for the tagset (see
2 Section 3.4), whose choice will have an important influence on the patterns that
3 are extracted, as well as for the settings (e.g. case sensitivity, mentioned also in
4 Section 3.4). There is no such thing as *the* n-grammar of English.

5.3 Providing a functional description

8 The quantitative measures discussed in Sections 5.1 and 5.2 may be interesting to
9 the numerically-minded learner. However, especially if we want to adopt a con-
10 structionist approach to language pedagogy, we should also attempt to show
11 how each pattern has its own functional properties, unless one is prepared to
12 argue for “the legitimacy of semantically null constructions” (Fillmore et al.
13 2012: 326; see Hilpert 2014: Chapter 3 for discussion). Providing a semantic or
14 functional characterization of a POS n-gram may need some inventiveness on
15 the part of the grammarian. Yet, by looking at the most frequent lexical n-gram
16 instantiations, we often get a clue as to what the pattern is predominantly used
17 for. In the case of the by now familiar pattern [*the* X_{noun} *of the* Y_{noun}], we could
18 formulate its function along the lines shown in (5):
19

20
21 (5) [*the* X_{noun} *of the* Y_{noun}]

22 This pattern allows speakers to link two entities: the noun phrase following
23 *of* (e.g. *the road*) and the noun preceding it (*side*). Among the most frequent
24 instantiations of this pattern, there are quite a few sequences where the
25 first noun denotes a portion (e.g. *rest*) or a position or dimension in space
26 or time (e.g. *end, side, edge, center, middle, back, top*) which ‘zooms in’
27 on a part of the larger whole expressed by the noun phrase after *of*. Not
28 surprisingly, this pattern overlaps with the next most productive one,
29 namely [*at/in/to/... the* X_{noun} *of the* Y_{noun}] (pattern No. 2), which adds a
30 preposition to indicate a relation to this spatial or temporal portion or
31 location, e.g. *at the end of the Y, in the middle of the Y, at the top of the Y,*
32 *by the end of the Y.*

33 Note that this pattern’s functional description is not only informed by its fre-
34 quent lexical instantiations; it also brings out the formal and functional related-
35 ness of this pattern with another one.
36

5.4 Providing opportunities for practice

38 Finally, let us offer some thoughts on how the selected grammar patterns can be
39 integrated into language learning activities aimed at consolidating the syntactic
40

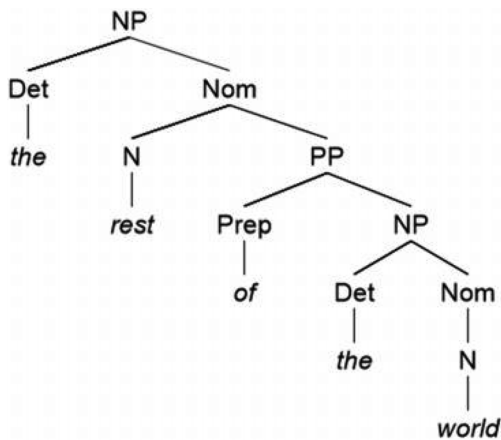
1 structures and their common instantiations. We hope that developers of lan-
2 guage learning materials might come up with a full range of concrete ideas, but
3 one obvious possibility of a practice activity for a given pattern is to encourage
4 learners to use suitable lexical instantiations in particular usage contexts. This
5 could be implemented as a simple fill-in exercise, which may or may not take
6 the form of a multiple-choice task whereby, given a particular sentence, learners
7 have to use the most suitable lexical n-gram from a set. These lexical n-grams
8 themselves, too, could be presented with gaps, which learners have to fill in
9 with contextually suitable items. Ideally, the sentences to be completed should
10 be taken from carefully selected authentic spoken or written discourse (although
11 one may have to clean up and/or simplify attested examples if learners are to
12 benefit from them optimally; cf. Gries 2008b); the fillers should be chosen from
13 the set of high-frequency sequences provided with the pattern (cf. Table 4).

14 Another sort of exercise could take the form of a role play between pairs of
15 students, who in a particular usage situation have to use a number of preselected
16 n-grams. For instance, two students could be asked to act as people of influence
17 in international politics, such as the Secretary of State of the US and the British
18 Prime Minister, discussing one or other rogue state's presumed possession of
19 weapons of mass destruction. Student A has to use (*give X / get / deserve*) *the*
20 *benefit of the doubt, the rest of the world, the State of the Union* and (*on/off/*
21 *from*) *the face of the earth*. Student B has to use *the end of the world, the fact of*
22 *the matter, (reach) the end of the line* and (*just*) *the tip of the iceberg*. The teacher
23 should not stop the role play until he or she is satisfied these sequences have
24 been used accurately and naturally (i.e. in a correct syntactic environment and,
25 whenever relevant, taking account of the idiomatic or encyclopedic meaning of
26 an expression, which of course should first have been illustrated by means of
27 authentic examples).

28 Once such common sequences have been mastered, a further exercise could
29 consist in using n-grams flexibly. For instance, teams could compete against
30 each other to produce the highest number of phrases that instantiate a POS
31 n-gram. This would allow them to learn the language by thinking in more
32 general categories, to exploit the combinatorial flexibility of grammar and use
33 and reuse at best available linguistic chunks. Other creatively-oriented exercises
34 could be to use patterns to form rhymes (for more advanced students), but of
35 course, such an exercise should not replace tasks that appeal more directly to
36 functional needs. As an alternative focusing again more on realistic language
37 use, students could be asked to detect instantiations of a set of POS n-grams,
38 say 5 different ones provided to them, in an authentic text. Some of these POS
39 n-grams, and accordingly their instantiations, may overlap (e.g. *with the tip of*
40 *the* and *the tip of the tongue*), demonstrating how n-grams can incrementally

1 combine to form full sentences. An easy related exercise could be to ask learners
 2 which lexical 5-grams in a text form complete constituents and which lexical
 3 5-grams do not. Such an exercise would raise students' awareness of language
 4 structure and might enhance linguistic insight.

5 For students of linguistics, the 'linear' approach to grammar proposed here
 6 could be offset by digressions on how the seemingly purely sequential structure
 7 of grammar patterns is actually hierarchically organized. An n-gram-based
 8 approach to grammar need not be incompatible with a more traditional linguistic
 9 reflection about constituent structure. For instance, we may explain to students
 10 that the sequence *the rest of the world* has the structure shown in Figure 7.



25 **Figure 7:** Tree diagram showing the hierarchical structure of a linear sequence

27 Such a tree diagram could lead to a discussion of recursion – the fact that a
 28 noun phrase can contain a noun phrase, for instance – or of why a preposition
 29 should not be simplistically defined as a word which comes before a noun –
 30 since in that case, a determiner would also be a preposition. An exercise could
 31 be to form complex trees by means of tree fragments constituting a toy grammar,
 32 such as one in which NP branches into Det and Nom, another in which Nom
 33 branches into N and PP (or just N), and yet another in which PP branches into
 34 P and NP.

36 6 Conclusion

39 These are exciting times. We have access to online corpora, specifically designed
 40 for linguistics or otherwise, containing staggering amounts of words. Corpora

1 have opened our eyes to lexical frequencies, collocations, and so on, but they
 2 need not close our eyes to syntactic patterns. With the help of automatic taggers
 3 that reach high precision rates, we can now hold our descriptions of grammar to
 4 the same empirical standard as our descriptions of the lexicon.

5 Espousing a radically usage-based approach to grammar, we have shown
 6 here how we can make use of relatively schematic templates derived from the
 7 Corpus of Contemporary American English as the basis of an ‘n-grammar’ of
 8 English. The selection of patterns to be included in such a grammar is based
 9 on the corpus frequency of part-of-speech 5-grams, which we consider to be
 10 constructions. We have demonstrated that a small number of high-frequency
 11 5-grams – just one hundred out of several tens of thousands – can cover quite
 12 a large portion of actually used 5-word strings (as well as strings of 2, 3 and 4
 13 words). This leads to a rather revolutionary approach to developing a pedagogical
 14 grammar: it breaks with traditional sequencing in grammars, which often deal
 15 first with everything related to the verb, then the noun, etc. – in this or another
 16 convenient (for largely conventional) order. Our motivation for following the
 17 order suggested by corpus frequencies is that it seems to us only common sense
 18 that those patterns that underlie the highest number of concrete word sequences
 19 should be presented before any others. While Leech (2011) may be right in taking
 20 a somewhat more considered stance regarding the principle “more frequent =
 21 more important to learn”, this is mainly because some researchers (cf. De Cock
 22 and Granger 2004) have noted that learners may overuse vocabulary items, such
 23 as *big* or *nice*. It is true that there are very common words which learners will
 24 soon come across, use themselves successfully and as a result start feeling
 25 rather too comfortable with. Leech’s reservation, however, applies less to more
 26 complex grammar patterns of the sort we have considered here, consisting of
 27 as many as five segments.

28 Our proposal to rank frequent and productive n-gram templates may help
 29 EFL/ESL material developers to select form-function patterns for active mastery.
 30 We hope to have illustrated or at least suggested how the construction of a book-
 31 length construct-i-con, discussing the hundred most frequent 5-gram templates,
 32 can lead to a fresh, empirically based and ultimately perhaps more relevant
 33 approach to teaching grammar.

34

35

36 Acknowledgements

37

38 Earlier versions of this paper were presented not just at the international con-
 39 ference Constructionist Approaches to Language Pedagogy (CALP, 8-9 November
 40 2013, Brussels, Belgium) but also at the University of Erlangen-Nurnberg,

1 Germany (1 July 2014) and the University of Tsukuba, Japan (8 July 2014). We
 2 thank the organizers of these events for giving us a stage, as well as members
 3 of the audience whose questions helped us formulate our ideas more clearly.
 4 We are also especially grateful to the editors of this volume and the anonymous
 5 reviewers for their constructive comments. All remaining inadequacies are ours
 6 alone.

9 References

- 11 Abott-Smith, Kirsten & Michael Tomasello. 2006. Exemplar-learning and schematization in a
 12 usage-based account of syntactic acquisition. *The Linguistic Review* 23. 275–290.
- 13 Baayen, R. Harald. 1989. *A corpus-based approach to morphological productivity: Statistical anal-*
 14 *ysis and psycholinguistic interpretation*. Ph.D. dissertation, Free University of Amsterdam.
- 15 Biber, Douglas, Susan Conrad & Viviana Cortes. 2004. *If you look at...: Lexical bundles in*
 16 *university teaching and textbooks*. *Applied Linguistics* 25(3). 371–405.
- 17 Biber, Douglas, Susan Conrad & Geoffrey Leech. 2002. *Longman student grammar of spoken*
 18 *and written English*. London: Longman.
- 19 Biber, Douglas, Stig Johansson, Geoffrey Leech, Susan Conrad & Edward Finegan. 1999. *Long-*
 20 *man grammar of spoken and written English*. London: Longman.
- 21 Boas, Hans C. 2003. *A constructional approach to resultatives*. Stanford: CSLI Publications.
- 22 Bybee, Joan L. 2006. From usage to grammar: The mind's response to repetition. *Language* 82
 23 (4). 711–733.
- 24 Cappelle, Bert. 2014. Review of Stefan Thim, *Phrasal verbs: The English verb-particle construction*
 25 *and its history*. Berlin & New York: Mouton de Gruyter. *English Language and Linguistics*
 26 18(3). 572–586.
- 27 Celce-Murcia, Marianne & Diane Larsen-Freeman. 1999. *The grammar book: An ESL/EFL teacher's*
 28 *course*. Boston, MA: Heinle & Heinle.
- 29 Chomsky, Noam. 1970. Remarks on nominalization. In Roderick A. Jacobs & Peter S. Rosenbaum
 30 (eds.), *Reading in English transformational grammar*, 184–221. Waltham: Ginn.
- 31 Croft, William. 1998. Linguistic evidence and mental representations. *Cognitive Linguistics* 9(2).
 32 151–173.
- 33 Croft, William & D. Allen Cruse. 2004. *Cognitive Linguistics*. Cambridge: Cambridge University
 34 Press.
- 35 Cowie, A.P. (ed.). 1998. *Phraseology: Theory, analysis, and applications*. Oxford: Oxford Univer-
 36 sity Press.
- 37 Davies, Mark. 2008–. *The Corpus of Contemporary American English: 450 million words, 1990–*
 38 *present*. Available online at <http://corpus.byu.edu/coca/>.
- 39 Davies, Mark. 2013. *Corpus of Global Web-Based English: 1.9 billion words from speakers in 20*
 40 *countries*. Available online at <http://corpus2.byu.edu/glowbe/>.
- 41 Davies, Mark & Dee Gardner. 2010. *A frequency dictionary of contemporary American English: Word sketches, collocates, and thematic lists*. London/New York: Routledge.
- 42 De Cock, Sylvie & Sylviane Granger. 2004. Computer learner corpora and monolingual learners' dictionaries: The perfect match. *Lexicographica* 20. 72–86.

- 1 De Kunder, Maurice. 2006. *Geschatte grootte van het geïndexeerde World Wide Web* [Estimated
2 size of the World Wide Web]. MA dissertation. Tilburg University.
- 3 Ellis, Nick C. 1996. Sequencing in SLA: Phonological memory, chunking, and points of order.
4 *Studies in Second Language Acquisition* 18. 91–126.
- 5 Ellis, Nick C. 2003. Constructions, chunking, and connectionism: The emergence of second lan-
6 guage structure. In Catherine J. Doughty & Michael H. Long (eds.), *Handbook of second
7 language acquisition*, 63–103. Malden, MA: Blackwell.
- 8 Ellis, Nick C. 2006. Cognitive perspectives on SLA: The Associative-Cognitive CREED. *AILA
9 Review* 19: 100–121.
- 10 Ellis, Nick C. 2013. Second language acquisition. In Thomas Hoffmann & Graeme Trousdale
11 (eds.), *Oxford handbook of Construction Grammar*, 365–378, Oxford: Oxford University
12 Press.
- 13 Ellis, Nick C. & Teresa Cadierno. 2009. Constructing a second language: Introduction to the
14 special section. *Annual Review of Cognitive Linguistics* 7. 111–139.
- 15 Fagan, Stephen & Ramazan Gençay. 2010. An introduction to textual econometrics. In Aman
16 Ullah & David E. A. Giles, *Handbook of empirical economics and finance*, 133–153. Boca
17 Raton: Chapman & Hall/CRC.
- 18 Fillmore, Charles J., Russell R. Lee-Goldman & Russell Rhodes. 2012. The FrameNet Construc-
19 tion. In Hans C. Boas & Ivan A. Sag (eds.), *Sign-Based Construction Grammar*, 283–299.
20 Stanford: CSLI.
- 21 Forsberg, Markus, Richard Johansson, Linnéa Bäckström, Lars Borin, Benjamin Lyngfelt, Joel
22 Olofsson & Julia Prentice. 2014. From construction candidates to constructicon entries:
23 An experiment using semi-automatic methods for identifying constructions in corpora.
24 *Constructions and Frames* 6(1). 114–135.
- 25 Francis, Gill, Susan Hunston & Elizabeth Manning. 1996. *Collins COBUILD Grammar Patterns 1:
26 Verbs*. London: HarperCollins. Available online at [http://arts-ccr-002.bham.ac.uk/ccr/
27 patgram/](http://arts-ccr-002.bham.ac.uk/ccr/patgram/).
- 28 Francis, Gill, Susan Hunston & Elizabeth Manning. 1998. *Collins COBUILD Grammar Patterns 2:
29 Nouns and adjectives*. London: HarperCollins.
- 30 Goldberg, Adele E. 1995. *Constructions: A Construction Grammar approach to argument struc-
31 ture*. Chicago: University of Chicago Press.
- 32 Goldberg, Adele E. 2003. Constructions: A new theoretical approach to language. *Trends in
33 Cognitive Sciences* 7(5). 219–224.
- 34 Goldberg, Adele E. 2006. *Constructions at work: The nature of generalization in language*.
35 Oxford: Oxford University Press.
- 36 Gries, Stefan Th. 2008a. Dispersions and adjusted frequencies in corpora. *International Journal
37 of Corpus Linguistics* 13(4). 403–437.
- 38 Gries, Stefan Th. 2008b. Corpus-based methods in analyses of second language acquisition
39 data. In Peter Robinson & Nick Ellis (eds.), *Handbook of Cognitive Linguistics and second
40 language acquisition*, 406–431. New York: Taylor & Francis.
- Guthrie, David, Ben Allison, Wei Liu, Louise Guthrie & Yorick Wilks. 2006. A closer look at skip-
gram modelling. *Proceedings of the fifth international conference on language resources
and evaluation (LREC'06)*, 1222–1225. Genoa, Italy.
- Halliday, M.A.K. 1978. *Language as social semiotic: The social interpretation of language and
meaning*. London: Arnold.
- Hilpert, Martin. 2014. *Construction Grammar and its application to English*. Edinburgh: Edin-
burgh University Press.

- 1 Huddleston, Rodney, Geoffrey K. Pullum et al. 2002. *The Cambridge grammar of the English*
2 *language*. Cambridge: Cambridge University Press.
- 3 Hunston, Susan & Gill Francis. 2000. *Pattern Grammar: A corpus-driven approach to the lexical*
4 *grammar of English*. Amsterdam: John Benjamins.
- 5 Iwata, Seizi. 2008. *Locative alternation: A lexical-constructional approach*. Amsterdam: John
6 Benjamins.
- 7 Jackendoff, Ray S. 1977. *X-bar syntax: A study of phrase structure*. Cambridge, MA: MIT Press.
- 8 Jurafsky, Dan. 1991. *An on-line computational model of human sentence interpretation: A theory*
9 *of the representation and use of linguistic knowledge*. Ph.D. Dissertation, University of
10 California, Berkeley.
- 11 Jurafsky, Dan. 2012. Language modeling: Introduction to n-grams. Lecture slides of Stanford
12 University's online course Natural Language Processing. [https://class.coursera.org/nlp/](https://class.coursera.org/nlp/lecture/14)
13 [lecture/14](https://class.coursera.org/nlp/lecture/14) (last accessed on 7 August 2014).
- 14 Kandel, Eric R. 2001. The molecular biology of memory storage: A dialogue between genes and
15 synapses. *Science* 294. 1030–1038.
- 16 Langacker, Ronald W. 2005. Construction grammars: Cognitive, radical, and less so. In Francisco
17 J. Ruiz de Mendoza Ibáñez & M. Sandra Peña Cervel (eds.), *Cognitive Linguistics: Internal*
18 *dynamics and interdisciplinary interaction*, 101–159. Berlin: Mouton de Gruyter.
- 19 Leech, Geoffrey. 2011. Frequency, corpora and language learning. In Fanny Meunier, Sylvie De
20 Cock, Gaëtanelle Gilquin & Magali Paquot (eds.), *A taste for corpora: In honour of Sylviane*
21 *Granger*, 7–32. Amsterdam: John Benjamins.
- 22 Levshina, Natalia & Kris Heylen. 2014. A radically data-driven Construction Grammar: Experi-
23 ments with Dutch causative constructions. In Ronny Boogaart, Timothy Coleman & Gijsbert
24 Rutten (eds.), *Extending the Scope of Construction Grammar*, 17–46. Berlin: Mouton de
25 Gruyter.
- 26 Lewis, Michael. 1993. *The Lexical Approach: The State of ELT and the Way Forward*. Hove, UK:
27 Language Teaching Publications.
- 28 Lin, Yuri, Jean-Baptiste Michel, Erez Lieberman Aiden, Jon Orwant, William Brockman & Slav
29 Petrov. 2012. Syntactic annotations for the Google Books Ngram Corpus. *Proceedings of*
30 *the 50th annual meeting of the Association for Computational Linguistics, Volume 2: Demo*
31 *papers (ACL '12)*, 169–174. Stroudsburg, PA: Association for Computational Linguistics.
- 32 Lyngfelt, Benjamin, Lars Borin, Markus Forsberg, Julia Prentice, Rudolf Rydstedt, Emma Sköldbog
33 & Sofia Tingsell. 2012. Adding a Constructicon to the Swedish resource network of
34 Språkbanken. *Proceedings of KONVENS 2012 (LexSem 2012 workshop)*, 452–461. Vienna,
35 September 2012.
- 36 McCarthy, Michael & Felicity O'Dell. 2001. *English vocabulary in use: Upper-intermediate*.
37 *Second Edition*. Cambridge: Cambridge University Press.
- 38 Michel, Jean-Baptiste, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew K. Gray,
39 William Brockman, The Google Books Team, Joseph P. Pickett, Dale Hoiberg, Dan Clancy,
40 Peter Norvig, Jon Orwant, Steven Pinker, Martin A. Nowak & Erez Lieberman Aiden. 2010.
Quantitative analysis of culture using millions of digitized books. *Science* 331(6014). 176–
182.
- O'Donnell, Matthew Brook. 2011. The adjusted frequency list: A method to produce cluster-
sensitive frequency lists. *ICAME Journal* 35. 135–169.
- Pulvermüller, Friedemann & Andreas Knoblauch. 2009. Discrete combinatorial circuits emerg-
ing in neural networks: A mechanism for rules of grammar in the human brain? *Neural*
Networks 22. 161–172.

- 1 Pulvermüller, Friedemann, Bert Cappelle & Yury Shtyrov. 2013. Brain basis of meaning, words,
2 constructions, and grammar. In Thomas Hoffmann & Graeme Trousdale (eds.), *Oxford*
3 *handbook of Construction Grammar*, 396–416, Oxford: Oxford University Press.
- 4 Robinson, Peter. 2006. Attention, memory, and the “noticing” hypothesis. *Language Learning*
5 45(2). 283–331.
- 6 Römer, Ute. 2010. Establishing the phraseological profile of a text type: The construction of
7 meaning in academic book reviews. *English Text Construction* 3(1). 95–119.
- 8 Schmidt, Richard W. 1990. The role of consciousness in second language learning. *Applied*
9 *Linguistics* 11(2). 129–158.
- 10 Sinclair, John. 1991. *Corpus, concordance, collocation: Describing English language*. Oxford:
11 Oxford University Press.
- 12 Tomasello, Michael. 2003. *Constructing a language: A usage-based theory of language acqui-*
13 *sition*. Cambridge, MA: Harvard University Press.
- 14 Wible, David & Tsao, Nai-Lung. 2010. StringNet as a computational resource for discovering and
15 investigating linguistic constructions. *Proceedings of the NAACL HLT workshop on extract-*
16 *ing and using constructions in computational linguistics*, 25–31. Los Angeles, CA: ACL.
- 17 Wray, Alison. 2002. *Formulaic language and the lexicon*. Cambridge: Cambridge University
18 Press.
- 19 Wulff, Stefanie & Stefan Th. Gries. 2011. Corpus-driven methods for assessing accuracy in
20 learner production. In Peter Robinson (ed.), *Second language task complexity: Research-*
21 *ing the cognition hypothesis of language learning and performance*, 61–88. Amsterdam:
22 Benjamins.
- 23 Zipf, George K. 1935. *The psychobiology of language: An introduction to dynamic philology*.
24 Cambridge, MA.: MIT Press.
- 25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40