

The adjusted frequency list: A method to produce cluster-sensitive frequency lists

Matthew Brook O'Donnell
University of Michigan

1 Introduction

The suggestion that language learners acquire and make use of multi-word chunks without either breaking them apart or building them up from individual words is well established in psycholinguistic (Pawley and Syder 1983; Ellis 1996, 2003) and corpus linguistic (Sinclair 1991; Stubbs 2002; Meunier and Granger 2008) circles. It is now even discussed in the popular press, as evidenced by a recent edition of the *New York Times* column *On Language* (Zimmer 2010). Frequency lists of items of various lengths are important in both computational and applied linguistics. They are also valuable for measuring the idiomatic/formulaic nature of text (Erman and Warren 2000; Sinclair and Mauranen 2006; Wray 2008). However, many of our computational tools and methods still focus on individual words as the foundational units of analysis.¹ The method proposed here is designed to address this issue.

In discussing the role of chunks in core vocabulary, particularly as it relates to language learners and language teaching, O'Keeffe, McCarthy and Carter highlight the fact that "many chunks are as frequent as or more frequent than the single-word items which appear in the core vocabulary" (2006: 46). Using the CANCODE corpus they found that only 33 single word items appear more frequently than the most common two-word chunk *you know*. Two of those 33 single-word items will be *you* and *know*, and for the latter it seems likely, as O'Keeffe *et al.* suggest, that its high ranking will be due in large part being a part of the highly frequent chunk in spoken English. Such observations highlight the importance of considering the role of chunks of two or more words in the description and teaching of vocabulary (also Nattinger and DeCarrico 1992).

A common methodological step in a corpus linguistic analysis is the extraction of frequency lists of various size chunks (variously called clusters, lexical bundles or n-grams). Most software packages facilitate the creation of such lists, making it possible to compare units of different length. However, each size unit

is (necessarily) counted on its own terms without reference to larger units of which they may be a part. For example every instance of *know* is counted individually even if each one of them is preceded by *you*, thus *you know* and *know* have the same frequency. This issue has been discussed with reference to larger units where collecting 3-, 4- and 5-grams together will result in very similar and often identical counts for *at the end*, *the end of*, *at the end of*, *the end of the*, *at the end of the*, *the end of the day* and so on.

The concept of the adjusted frequency list proposed here adjusts the frequency of items of various lengths when they are part of a larger unit that occurs at or above a given frequency or statistical threshold. That is, if *you know* occurs 15 times in a corpus and *know* 20 times, then the frequency of *know* will be adjusted from 20 down to five. The method outlined is ‘cluster sensitive’ because it boosts the rank of larger word sequences and builds on the notion that if such chunks are single choice items for speakers they should be counted as single items and their internal constituents left uncounted.

Section 2 provides a motivating example for the new procedure of counting words and n-grams which is described in Section 3. The next section describes three potential algorithms to implement the adjusted frequency list procedure. The second (using an index) and third (a two pass process) options are the better approaches and these are used in the case studies reported in Section 6. Two components of the BNC Baby corpus are examined by producing lists of 1- to 5-grams before and after the application of the adjusted frequency list procedure.

2 *First interlude: How does a corpus linguist tell a bedtime story?*

- (1) Once upon a time, there was a little girl named Goldilocks. She went for a walk in the forest. Pretty soon, she came upon a house. She knocked and, when no one answered, she walked right in...

Most readers will be familiar with how this text continues and recognize it as the story of “Goldilocks and the Three Bears” (see Text 1 in the Appendix for full text). How might a typical corpus linguist begin to ‘read’ (analyze) this particular text? Most likely he or she would begin by generating a word frequency list such as the one in Table 1. As is typical of just about any sample of English, the most frequent types are function words: *the*, *she*, *in*, *and*. These are followed by content words that give some key to who and what the story is about: *chair*, *porridge*, *bear*, *Goldilocks*. From this, therefore, we might answer that a corpus linguist would read this text one word at a time!

Table 1: Frequency list for Top 60 words from Text 1

<i>the</i>	34	<i>a</i>	5	<i>papa</i>	3
<i>she</i>	29	<i>just</i>	5	<i>ran</i>	3
<i>in</i>	14	<i>said</i>	5	<i>second</i>	3
<i>and</i>	13	<i>bears</i>	4	<i>sitting</i>	3
<i>chair</i>	10	<i>down</i>	4	<i>tasted</i>	3
<i>porridge</i>	10	<i>into</i>	4	<i>then</i>	3
<i>bear</i>	9	<i>of</i>	4	<i>there</i>	3
<i>been</i>	9	<i>right</i>	4	<i>they</i>	3
<i>my</i>	9	<i>sleeping</i>	4	<i>ahhh</i>	2
<i>someone's</i>	9	<i>up</i>	4	<i>as</i>	2
<i>too</i>	8	<i>all</i>	3	<i>ate</i>	2
<i>was</i>	8	<i>baby</i>	3	<i>bedroom</i>	2
<i>goldilocks</i>	7	<i>bowl</i>	3	<i>big</i>	2
<i>it</i>	7	<i>but</i>	3	<i>came</i>	2
<i>this</i>	7	<i>eating</i>	3	<i>cried</i>	2
<i>to</i>	7	<i>exclaimed</i>	3	<i>decided</i>	2
<i>bed</i>	6	<i>first</i>	3	<i>forest</i>	2
<i>is</i>	6	<i>growled</i>	3	<i>from</i>	2
<i>so</i>	6	<i>lay</i>	3	<i>home</i>	2
<i>three</i>	6	<i>mama</i>	3	<i>last</i>	2

However, a frequency list of single word items only tells part of the story. Like many stories written for and told to children, “Goldilocks and the Three Bears” contains a certain degree of repetition of phrases, for example, *the three bears*, *someone’s been eating my porridge*, *someone’s been sleeping in my bed*, *someone’s been sitting in my chair*, *growled the papa bear*, *said the mama bear*, *cried the baby bear*. The way to capture these kinds of chunks is to generate a frequency list of n-grams or clusters. While this is often done producing different lists for different values of *n*, it can be valuable to produce a single list covering a range of *n* values. This allows for the kind of comparison between single words and larger chunks alluded to in the quote above from O’Keeffe *et al.* (2006). Table 2 shows such a combined list of 1-, 2- and 3-grams for our bedtime story.² The top ten items in the list are still single words but 22 (37%) of the top 60 types are now clusters of two or three words. This suggests the impor-

tance of clusters in this text. Now we might want to answer that a corpus linguist would read the text in words AND chunks at the same time.

Table 2: Frequency list of Top 60 1-, 2- and 3-grams from Text 1

<i>the</i>	34	<i>is</i>	6	<i>baby</i>	3
<i>she</i>	29	<i>so</i>	6	<i>baby bear</i>	3
<i>in</i>	14	<i>so she</i>	6	<i>been eating</i>	3
<i>and</i>	13	<i>three</i>	6	<i>been eating my</i>	3
<i>chair</i>	10	<i>a</i>	5	<i>been sitting</i>	3
<i>porridge</i>	10	<i>just</i>	5	<i>been sitting in</i>	3
<i>bear</i>	9	<i>said</i>	5	<i>been sleeping</i>	3
<i>been</i>	9	<i>bears</i>	4	<i>been sleeping in</i>	3
<i>my</i>	9	<i>down</i>	4	<i>bowl</i>	3
<i>someone's</i>	9	<i>into</i>	4	<i>but</i>	3
<i>someone's been</i>	9	<i>is too</i>	4	<i>chair is</i>	3
<i>too</i>	8	<i>of</i>	4	<i>eating</i>	3
<i>was</i>	8	<i>right</i>	4	<i>eating my</i>	3
<i>goldilocks</i>	7	<i>sleeping</i>	4	<i>eating my porridge</i>	3
<i>in the</i>	7	<i>the three</i>	4	<i>exclaimed</i>	3
<i>it</i>	7	<i>the three bears</i>	4	<i>first</i>	3
<i>this</i>	7	<i>three bears</i>	4	<i>growled</i>	3
<i>to</i>	7	<i>up</i>	4	<i>in my bed</i>	3
<i>bed</i>	6	<i>all</i>	3	<i>in my chair</i>	3
<i>in my</i>	6	<i>and she</i>	3	<i>into the</i>	3

But again this answer is not without some limitations. Notice how the words *been* and *someone's* have the same frequency (9 occurrences) individually as the bigram *someone's been*. Similarly *eating*, *eating my* and *eating my porridge* all have a frequency of three and likewise *bears*, *the three*, *three bears* and *the three bears*. There are six occurrences of both *so* and *so she* and three of both *baby* and *baby bear*. In each of these cases the largest n-gram accounts for all the occurrences of the smaller n-grams and single words. This raises the question of whether the smaller units should really be included in the frequency list or not. In other instances most but not all of the occurrences of a word can be accounted for by a larger cluster. For example, *into* occurs four times in Table 2 and *into*

the three times, leaving just one instance of *into* not accounted for by the bigram. The same goes for *been sleeping in* (3 occurrences) and *sleeping* (4 occurrences). In these instances the individual words should certainly remain in the frequency list but their rank appears to be inflated because of the larger cluster.

So is our intrepid corpus linguist perhaps over reading (analyzing) the individual words in the story? How might this issue be addressed?

3 A new concept for frequency counts: The adjusted frequency list

On the wall of my office I have a Dr Seuss ABC poster similar to those often found in a child's nursery or toddler's bedroom. It reads: *A is for Alligator, B is for Ball, C is for Cat*, and so on. Consider the following 'text' (Text 2), which is 14 tokens long constructed using the first five types:

- (2) *Alligator Ball Cat Alligator Ball Cat Alligator Ball Duck Alligator Elephant Alligator Ball Cat*

Table 3 contains the frequency lists for all the 1-, 2- and 3- grams in this text. The lists are ordered by frequency and then alphabetically.

Table 3: Frequency lists of all 1-, 2- and 3-grams from Text 2

Words (1-grams)		2-grams	
<i>Alligator</i>	5	<i>Alligator Ball</i>	4
<i>Ball</i>	4	<i>Ball Cat</i>	3
<i>Cat</i>	3	<i>Cat Alligator</i>	2
<i>Duck</i>	1	<i>Alligator Elephant</i>	1
<i>Elephant</i>	1	<i>Ball Duck</i>	1
		<i>Duck Alligator</i>	1
		<i>Elephant Alligator</i>	1
3-grams			
<i>Alligator Ball Cat</i>		3	
<i>Ball Cat Alligator</i>		2	
<i>Cat Alligator Ball</i>		2	
<i>Alligator Ball Duck</i>		1	
<i>Alligator Elephant Alligator</i>		1	
<i>Ball Duck Alligator</i>		1	
<i>Duck Alligator Elephant</i>		1	
<i>Elephant Alligator Ball</i>		1	

If these three lists are merged (again on the basis of frequency and then alphabetically) the list in Table 4 results. For this text the most frequent bigram (*Alligator Ball*) shares the same rank as the second most frequent single item (*Ball*). Similarly, *Alligator Ball Cat*, the most frequent trigram has the same frequency as the second most frequent bigram (*Ball Cat*) and third most frequent single item (*Cat*). From a vocabulary analysis perspective this reinforces the point made by O’Keeffe *et al.* (2006) regarding the value of including clusters in banded frequency lists. N-gram lists are built using a moving window of one word at a time through the text and counting units of length *n*, e.g. *Alligator Ball*, *Ball Cat*, *Cat Alligator* (with *n=2*). This means that aside from the first and last word of a text when collecting units of length *n*, each word is counted *n* times.

Table 4: Combined Frequency list of all 1-, 2- and 3-grams from Text 2

<i>Alligator</i>	5
<i>Alligator Ball</i>	4
<i>Ball</i>	4
<i>Alligator Ball Cat</i>	3
<i>Ball Cat</i>	3
<i>Cat</i>	3
<i>Ball Cat Alligator</i>	2
<i>Cat Alligator</i>	2
<i>Cat Alligator Ball</i>	2
<i>Alligator Ball Duck</i>	1
<i>Alligator Elephant</i>	1
<i>Alligator Elephant Alligator</i>	1
<i>Ball Duck</i>	1
<i>Ball Duck Alligator</i>	1
<i>Duck</i>	1
<i>Duck Alligator</i>	1
<i>Duck Alligator Elephant</i>	1
<i>Elephant</i>	1
<i>Elephant Alligator</i>	1
<i>Elephant Alligator Ball</i>	1

One of the uses of an n-gram list is to discover recurring units that might be formulaic or idiomatic and function as a single choice for the language user (cf. the ‘idiom choice principle’, Sinclair 1991; Erman and Warren 2000). Setting a threshold for recurrence balances the over counting of the moving window procedure and also serves as a crude measure of formulaicity.

In order to simplify things, consider for a moment a frequency list of Text 2 with all the single words and just the bigrams occurring at least three times (see Table 5).

Table 5: Combined Frequency list of all words and the 2-grams with frequency > 2 in Text 2

<i>Alligator</i>	5
<i>Alligator Ball</i>	4
<i>Ball</i>	4
<i>Ball Cat</i>	3
<i>Cat</i>	3
<i>Duck</i>	1
<i>Elephant</i>	1

What this list suggests is that the bigrams *Alligator Ball* and *Ball Cat* are actually single choice units. Ignoring the fact that there is overlap between the units (*Alligator Ball* always overlaps with *Ball Cat*) the text becomes:

(2b) *Alligator Ball* *Ball Cat* *Alligator Ball* *Ball Cat* *Alligator Ball*
 Duck Alligator Elephant *Alligator Ball* Cat

(where *Alligator Ball* and *Ball Cat* indicate single units). With the text viewed in this manner the resulting frequency list, shown in Table 6, contains five types and ten tokens:

Table 6: Adjusted Frequency list of all words and the 2-grams with frequency > 2 in Text 2b

<i>Alligator Ball</i>	4
<i>Ball Cat</i>	3
<i>Alligator</i>	1
<i>Duck</i>	1
<i>Elephant</i>	1

Notice that single items *Ball* and *Cat* have disappeared from the list because they do not appear independently of the clusters *Alligator Ball* and *Ball Cat*. The count for *Alligator* is reduced from five to just a single occurrence because of the four instances of *Alligator Ball*. I propose the term ‘adjusted frequency list’ for a frequency list that has undergone this kind of adjustment.

Now what happens if we include trigrams into consideration while keeping the same threshold of three or more occurrences for n-grams. This adds only one item, *Alligator Ball Cat*, to the unadjusted frequency list (see Table 7):

Table 7: Combined Frequency list of all words and the 2- and 3-grams with frequency > 2 in Text 2

<i>Alligator</i>	5
<i>Alligator Ball</i>	4
<i>Ball</i>	4
<i>Alligator Ball Cat</i>	3
<i>Ball Cat</i>	3
<i>Cat</i>	3
<i>Duck</i>	1
<i>Elephant</i>	1

Now applying the same adjustment procedure in which longer units (*Alligator Ball Cat*) should take precedence over their component parts (*Alligator Ball*, *Ball Cat*, *Alligator*, *Ball* and *Cat*), the text now consists of five types and seven tokens.

(2c) *Alligator Ball Cat* *Alligator Ball Cat* *Alligator Ball*
 Duck Alligator Elephant *Alligator Ball Cat*

Table 8 contains the adjusted frequency list for Text 2 using a frequency threshold of 3 occurrences for n-grams (with $n > 1$). As before single items *Ball* and *Cat* have disappeared and *Alligator* is reduced to a single occurrence. The bigram *Ball Cat* has been removed because it does not occur independently of the trigram *Alligator Ball Cat*, and the four occurrences of *Alligator Ball* have been reduced to the single instance where the bigram is not followed by *Cat*.

Table 8: Adjusted Frequency list of all words and the 2- and 3-grams with frequency > 2 in Text 2c

<i>Alligator Ball Cat</i>	3
<i>Alligator</i>	1
<i>Alligator Ball</i>	1
<i>Duck</i>	1
<i>Elephant</i>	1

Although only a toy example, it should be sufficient to illustrate the notion of the adjusted frequency list. There are a number of parameters, particularly the thresholds to use at various values of n and the maximum value of n , that can be tuned and will result in different outputs. But the key characteristic of the procedure is that it is sensitive to the use of clusters as (potentially) single lexical choices.

4 Algorithms for the adjusted frequency list procedure

The previous section provided an overview of the adjusted frequency list procedure without any suggestion of how it might be implemented. This section presents three possible algorithms in some detail. It is not necessary to follow through the details of each algorithm and this whole section can be skimmed over without losing the overall concept of the adjusted frequency list.

4.1 Simple non-indexed algorithm

The first and simplest approach is applied just to the frequency list of 1-, 2-, ... N_{\max} -grams. Given a text or set of texts an adjusted frequency list is constructed in the following manner.

1. Construct frequency lists (or a single combined list) for all items length 1 to N_{\max} using the standard moving word window method and no frequency threshold (i.e. all items down to single occurrence).
2. Remove all items of length 2 to N_{\max} that occur with frequency less than desired threshold adopted for formula/unit status.
3. For each remaining n-gram with frequency f (in descending order by length, i.e. N_{\max} to 2) derive each of its component sub-items.
So for the trigram *Alligator Ball Cat* there are bigrams *Alligator Ball* and *Ball Cat* and three single items *Alligator*, *Ball* and *Cat*.
4. Reduce the frequency of each of these sub-items by f (unless frequency=0).

In essence this algorithm groups all items in a combined frequency list into a tree (or directed graph) with larger n-grams higher up the tree linked to smaller n-grams that are component parts of the larger n-gram. But the trigram *Alligator Ball Cat* will link not only to bigrams *Alligator Ball* and *Ball Cat* but also to each of individual words *Alligator*, *Ball* and *Cat*.

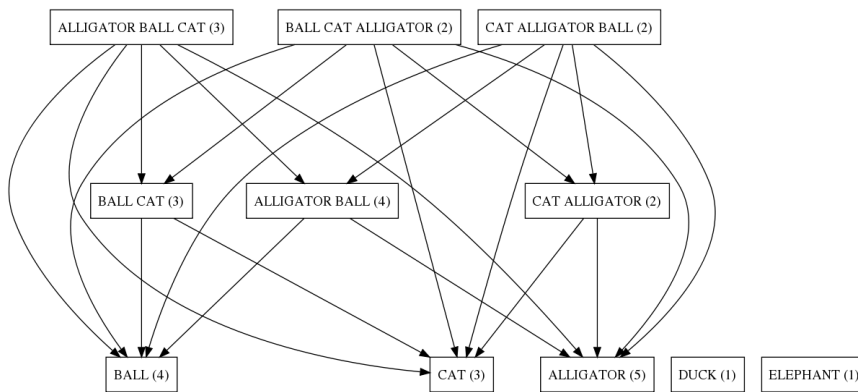
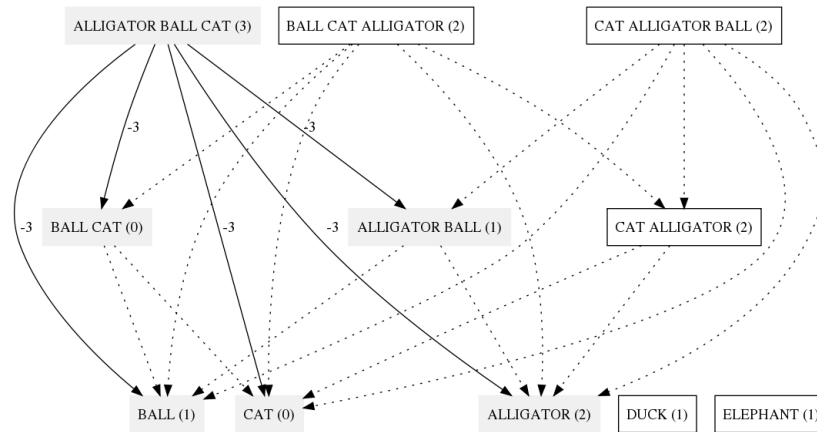


Figure 1: Links between n-grams in simple non-indexed algorithm

Figure 1 illustrates these connections for the 2- and 3-grams in Text 2 with a frequency of 2 or more extracted from Table 4. Figure 2 shows the first two iterations of the algorithm at Step 3 for the n-grams *Alligator Ball Cat* and *Ball Cat Alligator*.

Reduction process (Step 3)

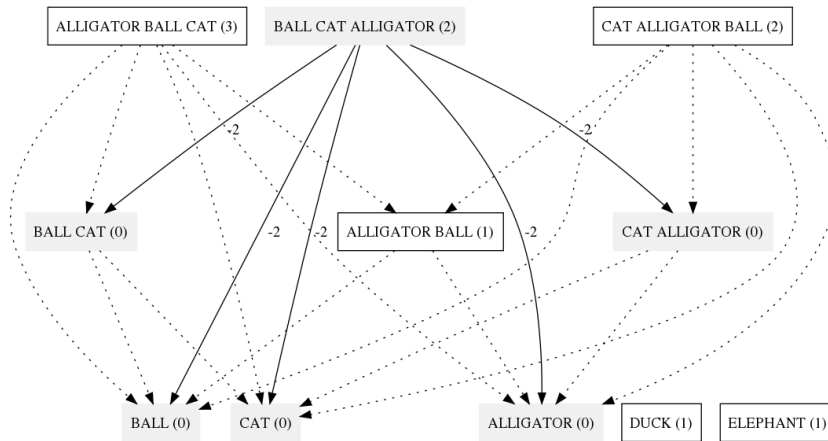


Resulting frequency list

- 3 Alligator Ball Cat**
- 2 Alligator
- 2 Ball Cat Alligator
- 2 Cat Alligator
- 2 Cat Alligator Ball
- 1 Alligator Ball
- 1 Ball
- 1 Duck
- 1 Elephant
- 0 Ball Cat
- 0 Cat

Figure 2a: Applying the simple non-indexed algorithm to n-grams in Text 2.

Reduction process (Step 3)



Resulting frequency list

- 3 *Alligator Ball Cat*
- 2 *Ball Cat Alligator***
- 2 *Cat Alligator Ball*
- 1 *Alligator Ball*
- 1 *Duck*
- 1 *Elephant*
- 0 *Alligator*
- 0 *Ball*
- 0 *Ball Cat*
- 0 *Cat*
- 0 *Cat Alligator*

Figure 2b: Applying the simple non-indexed algorithm to n-grams in Text 2

The problem with this simple method is that it is likely to be too productive in the final step. That is, given trigrams $[Alligator\ Ball_i\ Cat]$, $[Ball_i\ Cat\ Duck]$, $[Alligator\ Ball_j\ Cat]$ and $[Ball_j\ Cat\ Duck]$, generated from the string *Alligator*

$Ball_i$ *Cat Duck Alligator* $Ball_j$ *Cat Duck*, the counts for both $Ball_i$ and $Ball_j$ will be reduced twice. This is because the procedure has no knowledge of which particular *Ball* is being referenced. This is further illustrated in Figure 2, where after applying reductions to just two trigrams (*Alligator Ball Cat* and *Ball Cat Alligator*) the count for the single item *Alligator* has been reduced to zero. We know from Table 8 that the final count for *Alligator* should actually be one after applying the full procedure.

4.2 Indexed algorithm

To address the limitation of the simplest possible algorithm two further algorithms are presented. The first builds an index from the corpus and then can selectively reduce counts for smaller values of n as it reduces a specific n -gram. Given a text or set of texts an adjusted frequency list is constructed in the following manner:

1. Construct indexed frequency lists for all items length 1 to N_{\max} , so that each instance of an item is recorded with reference to its source file and position within that file (either just start or both start and end offsets).
2. Remove all items of length 2 to N_{\max} that occur less than desired threshold used for formula/unit status.
3. For each remaining n -gram with frequency f (in descending order by length, i.e. N_{\max} to 2) derive each of its component sub-items, recording the start and end positions for each occurrence of the n -gram.
4. For each of the sub-items identified in Step 3, scan their index records for an occurrence that falls within the position range of the larger n -gram and remove record.

So given Text 3:

(3) *Alligator*₁ *Ball*₂ *Cat*₃ *Alligator*₄ *Ball*₅ *Cat*₆ *Alligator*₇ *Ball*₈ *Cat*₉

where the subscripts indicate word (or position) offset, the following index entries would result from Step 1. (Each instance of an item is recorded with the form startOffset:endOffset).

<i>Alligator Ball Cat</i>	[1:3, 4:6, 7:9]	3
<i>Ball Cat Alligator</i>	[2:4, 5:7]	2
<i>Cat Alligator Ball</i>	[3:5, 6:8]	2
<i>Alligator Ball</i>	[1:2, 4:5, 7:8]	3
<i>Ball Cat</i>	[3:4, 5:6, 8:9]	3
<i>Cat Alligator</i>	[3:4, 6:7]	2
<i>Alligator</i>	[1:1, 4:4, 7:7]	3
<i>Ball</i>	[2:2, 5:5, 8:8]	3
<i>Cat</i>	[3:3, 6:6, 9:9]	3

After the application of Step 2 with a threshold of 3, the index would be:

<i>Alligator Ball Cat</i>	[1:3, 4:6, 7:9]	3
<i>Alligator Ball</i>	[1:2, 4:5, 7:8]	3
<i>Ball Cat</i>	[3:4, 5:6, 8:9]	3
<i>Alligator</i>	[1:1, 4:4, 7:7]	3
<i>Ball</i>	[2:2, 5:5, 8:8]	3
<i>Cat</i>	[3:3, 6:6, 9:9]	3

Step 3 would begin with the trigram *Alligator Ball Cat* and derive the sub-items *Alligator Ball*, *Ball Cat*, *Alligator*, *Ball* and *Cat*. For each entry in the index for *Alligator Ball Cat* the entries for these sub-items is scanned for entries that fall within the start and end offsets. Matching entries are deleted, as follows:

<i>Alligator Ball Cat</i>	[1:3 , 4:6, 7:9]	3
<i>Alligator Ball</i>	[1:2 , 4:5, 7:8]	2
<i>Ball Cat</i>	[2:3 , 5:6, 8:9]	2
<i>Alligator</i>	[1:1 , 4:4, 7:7]	2
<i>Ball</i>	[2:2 , 5:5, 8:8]	2
<i>Cat</i>	[3:3 , 6:6, 9:9]	2

<i>Alligator Ball Cat</i>	[1:3, 4:6 , 7:9]	3
<i>Alligator Ball</i>	[4:5 , 7:8]	1
<i>Ball Cat</i>	[5:6 , 8:9]	1
<i>Alligator</i>	[4:4 , 7:7]	1
<i>Ball</i>	[5:5 , 8:8]	1
<i>Cat</i>	[6:6 , 9:9]	1

<i>Alligator Ball Cat</i>	[1:3, 4:6, 7:9]	3
<i>Alligator Ball</i>	[7:8]	0
<i>Ball Cat</i>	[8:9]	0
<i>Alligator</i>	[7:7]	0
<i>Ball</i>	[8:8]	0
<i>Cat</i>	[9:9]	0

In this extreme case the adjusted frequency list contains a single trigram *Alligator Ball Cat* with a frequency of 3. All instances of the 5 sub-items, *Alligator Ball*, *Ball Cat*, *Alligator*, *Ball* and *Cat* occur within these the three instances of the trigram.

4.3 The Serial Cascading Algorithm

An alternative approach that does not need an index but avoids the problems of the non-indexed approach discussed in Section 4.1 has been suggested by Catherine Smith (p.c.). It takes two passes over the texts in a corpus. The first pass constructs the relevant n-gram lists and the second pass counts n-grams according to a largest n first cascade:

Pass #1

1. Construct frequency lists (or a single combined list) for all items length 2 to N_{\max} using the standard moving word window method and no frequency threshold (i.e. all items down to single occurrence).
2. Remove all items of length 2 to N_{\max} that occur less than desired threshold used for formula/unit status.

Pass #2

3. Initialize:
 - a. $\text{adjusted_list} = \{\}$
 - b. $p = 1$
 - c. $\text{last}_i = \emptyset$
4. Step through using a moving window of one token steps using position counter p .
5. Select $n\text{-gram}_{\text{candidate}}$, an n -gram of N_{max}
6. Check whether $n\text{-gram}_{\text{candidate}}$ is found in lists constructed in PASS #1
 - a. If yes and $p + N_{\text{max}} - 1 > \text{last}_i$ add one to the count for $n\text{-gram}_{\text{candidate}}$ in the adjusted list, set last_i to $p + N_{\text{max}}$ and return to Step 4
 - b. else reduce N_{max} by 1
 - i. If $N_{\text{max}} > 1$
 1. If $p + N_{\text{max}} - 1 > \text{last}_i$ return to Step 5
 2. else reset N_{max} and return to Step 4
 - ii. else if $p > \text{last}_i$ add one to single word count in the adjusted list and return to Step 4

If this algorithm is applied to Text 2 (used in Section 3) using a frequency threshold of three or greater for 2- and 3-grams the algorithm proceeds as follows:

Pass #1

Collect all 2- and 3-grams occurring three or more times in text.

<i>Alligator Ball</i>	4
<i>Alligator Ball Cat</i>	3
<i>Ball Cat</i>	3

Pass #2

The second part of the algorithm is somewhat complex. Below three snapshots of the process as applied to Text 2 are illustrated with the value of variables at each step shown:

Step	Location in text (n-gram _{candidate} in bold)	Variables	Adjusted list
3	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> <i>Alligator Ball Duck</i> <i>Alligator Elephant</i> <i>Alligator Ball Cat</i>	$p=1$ $last_i = \emptyset$	{ }
4 5 6 6a	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> <i>Alligator Ball Duck</i> <i>Alligator Elephant</i> <i>Alligator Ball Cat</i>	$p=1$ $last_i = \emptyset$ $N_{max}=3$ <i>Alligator Ball Cat</i> on list $last_i=3$	{ 'Alligator Ball Cat': 1 }
4 5 6 6b 6b.i. 2	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> <i>Alligator Ball Duck</i> <i>Alligator Elephant</i> <i>Alligator Ball Cat</i>	$p=2$ $last_i=3$ $N_{max}=3$ <i>Ball Cat Alligator</i> not on list $N_{max}=2$ $p + N_{max} - 1 = last_i$	{ 'Alligator Ball Cat': 1 }
4 5 6 6b 6b.i. 1	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> <i>Alligator Ball Duck</i> <i>Alligator Elephant</i> <i>Alligator Ball Cat</i>	$p=3$ $last_i=3$ $N_{max}=3$ <i>Cat Alligator Ball</i> not on list $N_{max}=2$ $p + N_{max} - 1 > last_i$	{ 'Alligator Ball Cat': 1 }
5 6 6b 6b.i. 6b.i. 2	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> <i>Alligator Ball Duck</i> <i>Alligator Elephant</i> <i>Alligator Ball Cat</i>	$p=3$ $last_i=3$ $N_{max}=2$ <i>Cat Alligator</i> not on list $N_{max}=1$ $p + N_{max} - 1 = last_i$ $N_{max}=3$	{ 'Alligator Ball Cat': 1 }
4 5 6 6a	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> <i>Alligator Ball Duck</i> <i>Alligator Elephant</i> <i>Alligator Ball Cat</i>	$p=4$ $last_i=3$ $N_{max}=3$ <i>Alligator Ball Cat</i> on list $last_i=6$	{ 'Alligator Ball Cat': 1 }
<i>(some intervening steps skipped)</i>			

5 6 6a	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> Alligator Ball Duck <i>Alligator Elephant</i> <i>Alligator Ball Cat</i>	p=7 last _i =6 N _{max} =2 Alligator Ball on list p + 2 - 1 > last _i last _i =8 N _{max} =3	{ 'Alligator Ball Cat': 2, 'Alligator Ball': 1 }
4 5 6 6b 6b.i 1	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> Alligator Ball Duck <i>Alligator Elephant</i> <i>Alligator Ball Cat</i>	p=8 last _i =8 N _{max} =3 Ball Duck Alligator not on list N _{max} =2 p + 2 - 1 > last _i	{ 'Alligator Ball Cat': 2, 'Alligator Ball': 1 }
5 6 6b 6b.i 6b.i 2	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> Alligator Ball Duck <i>Alligator Elephant</i> <i>Alligator Ball Cat</i>	p=8 last _i =8 N _{max} = 2 Ball Duck not on list N _{max} =1 p + 1 - 1 = last _i N _{max} =3	{ 'Alligator Ball Cat': 2, 'Alligator Ball': 1 }
(some intervening steps skipped)			
5 6 6b 6b.ii	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> <i>Alligator Ball Duck</i> Alligator Elephant <i>Alligator Ball Cat</i>	p=11 last _i =10 N _{max} = 2 Elephant Alligator not on list N _{max} =1 p > last _i	{ 'Alligator Ball Cat': 2, 'Alligator Ball': 1, 'Duck': 1, 'Alligator': 1, 'Elephant': 1 }
4 5 6 6a	<i>Alligator Ball Cat</i> <i>Alligator Ball Cat</i> <i>Alligator Ball Duck</i> <i>Alligator Elephant</i> Alligator Ball Cat	p=12 last _i =10 N _{max} = 3 Alligator Ball Cat on list last _i =11	{ 'Alligator Ball Cat': 3, 'Alligator Ball': 1, 'Duck': 1, 'Alligator': 1, 'Elephant': 1 }

The resulting output is:

<i>Alligator Ball Cat</i>	3
<i>Alligator Ball</i>	1
<i>Alligator</i>	1
<i>Duck</i>	1
<i>Elephant</i>	1

Future work is needed to implement and test these (and other algorithms) to gauge the applications for which each is best suited. The use of an index requires more in terms of computation resources but does allow for comparative concordancing of unadjusted and adjusted items. The Serial Cascading Algorithm is more lightweight and could potentially scale to a distributed/parallel implementation.

5 *Second interlude: The well-adjusted bedtime story*

We now return to our corpus linguist endeavoring to tell a bedtime story using the state-of-the-art tools of the trade. When we left them back in Section 2, they had begun to come to terms with the highly repetitious and chunky nature of the typical bedtime story and created a combined 1- to 3-gram list (see Table 2). But there were at least two problems with this approach. First, single words still fill the top ranks of the list even though many of them are components of highly frequent chunks of two or three words. And second, a number of frequent bigrams on the list were entirely accounted for by certain trigrams.

Table 9 shows the top of the adjusted frequency list for Text 1 for 1-, 2- and 3-grams using a frequency threshold of three or more occurrences for the inclusion of 2- and 3-grams in the adjustment process. When compared to the unadjusted list in Table 2 notice the marked reduction for *the* from 34 occurrences down to nine. This indicates that 25 instances of *the* are a part of bi- or trigram that occurs three times or more. Likewise the 9 occurrences of *my* (rank 9 in Table 1) are all accounted for by the trigrams: *eating my porridge* (3), *in my chair* (3) and *in my bed* (3). Notice also how the bigrams *been eating*, *been sitting*, *been sleeping*, all with three occurrences in Table 2, no longer occur in the adjusted list. This is because they are fully accounted for by the larger trigrams *been eating my*, *been sitting in* and *been sleeping in*.

Table 9: Adjusted Frequency list of Top 60 1-, 2- and 3-grams from Text 1

<i>she</i>	16	<i>been eating my</i>	3	<i>sleeping in my</i>	3
<i>and</i>	10	<i>been sitting in</i>	3	<i>someone's been eating</i>	3
<i>the</i>	9	<i>been sleeping in</i>	3	<i>someone's been sitting</i>	3
<i>goldilocks</i>	7	<i>bowl</i>	3	<i>someone's been sleeping</i>	3
<i>in the</i>	7	<i>but</i>	3	<i>the first</i>	3
<i>so she</i>	6	<i>eating my porridge</i>	3	<i>the mama bear</i>	3
<i>a</i>	5	<i>exclaimed</i>	3	<i>the second</i>	3
<i>was</i>	5	<i>growled</i>	3	<i>then</i>	3
<i>chair</i>	4	<i>in my bed</i>	3	<i>there</i>	3
<i>down</i>	4	<i>in my chair</i>	3	<i>they</i>	3
<i>is too</i>	4	<i>into the</i>	3	<i>this chair is</i>	3
<i>of</i>	4	<i>it all</i>	3	<i>this porridge is</i>	3
<i>porridge</i>	4	<i>it was</i>	3	<i>to the</i>	3
<i>the three bears</i>	4	<i>just right</i>	3	<i>a little</i>	2
<i>to</i>	4	<i>papa bear</i>	3	<i>ahhh</i>	2
<i>too</i>	4	<i>ran</i>	3	<i>ahhh this</i>	2
<i>up</i>	4	<i>said the mama</i>	3	<i>all up</i>	2
<i>and she</i>	3	<i>she lay</i>	3	<i>and ran</i>	2
<i>baby bear</i>	3	<i>she tasted the</i>	3	<i>and when</i>	2
<i>bed</i>	3	<i>sitting in my</i>	3	<i>as</i>	2

The five instances of *just* in Table 1 become three of *just right* and two for single-item *just* in Table 9. A side-effect of this grouping is that the two adverbial usages of *just* have been distinguished: i. exactly (*just right*) in lines 1,2 and 4 and ii. temporal in lines 3 and 5.

1 "Ahhh, this porridge is **just right**," she said happily and
 2 ir."Ahhh, this chair is **just right**," she sighed. But just
 3 right," she sighed. But **just** as she settled down into the
 4 he third bed and it was **just right**. Goldilocks fell aslee
 5 !" exclaimed Baby bear. **Just** then, Goldilocks woke up an

It is worth calling attention to a couple of points arising from the adjusted list in Table 9 that illustrate the effects of choices made with regards to the largest n-

gram (N_{\max}) included in the adjustment procedure and also the threshold (or thresholds) chosen for the different values of n . Here the procedure was applied at $N_{\max}=3$ with a threshold of 3+ occurrences. As a result there remains some overlapping n -grams that are actually part of a larger chunk. For instance, *been eating my* and *someone's been eating* both have 3 occurrences in the list in Table 9. These are clearly part of a larger 4-gram *someone's been eating my*. Likewise towards the end of the list we can see overlap between words and bigrams with frequencies below the selected threshold—for example: *ahh* and *ahh this* with 2 occurrences.³

These minor caveats aside, our corpus linguist now has a tool that provides a more realistic picture of the interaction of chunks and single words in the Goldilocks text. And the example of *just* demonstrates the potential of improved efficiency in a KWIC analysis, which as everyone knows is both the next act in the story and another story all by itself (see O'Donnell 2008).

6 Looking at some larger corpora

The final two examples apply the adjusted frequency list method to two of the categories in the BNC Baby sample corpus. These two sections are the 1 million word demographically sampled spoken component (30 texts) and the 1 million word sub-corpus of academic texts (also 30 texts).

6.1 BNC Baby Demographic section

The list in Table 10 contains the top 150 1-, 2- and 3-grams from the Demographic section of the corpus with no adjustment.

Table 10: Top 150 combined 1-, 2- and 3-grams according to type frequency in BNC Baby Demographic section

Rank	Item	Freq.	Rank	Item	Freq.	Rank	Item	Freq.
1	<i>i</i>	30371	51	<i>up</i>	4056	101	<i>it was</i>	1901
2	<i>you</i>	29688	52	<i>with</i>	3833	102	<i>very</i>	1878
3	<i>the</i>	27698	53	<i>erm</i>	3813	103	<i>can't</i>	1868
4	<i>it</i>	21834	54	<i>them</i>	3670	104	<i>five</i>	1855
5	<i>and</i>	19845	55	<i>at</i>	3662	105	<i>four</i>	1820
6	<i>a</i>	19600	56	<i>are</i>	3652	106	<i>on the</i>	1789
7	<i>to</i>	17180	57	<i>me</i>	3607	107	<i>been</i>	1770
8	<i>that</i>	14722	58	<i>you know</i>	3605	108	<i>bit</i>	1715

9	<i>yeah</i>	14303	59	<i>said</i>	3563	109	<i>alright</i>	1703
10	<i>oh</i>	10398	60	<i>two</i>	3528	110	<i>would</i>	1657
11	<i>in</i>	10133	61	<i>your</i>	3448	111	<i>him</i>	1655
12	<i>no</i>	9804	62	<i>out</i>	3168	112	<i>they're</i>	1653
13	<i>of</i>	9799	63	<i>i'm</i>	3153	113	<i>were</i>	1625
14	<i>it's</i>	8534	64	<i>see</i>	3143	114	<i>i know</i>	1623
15	<i>well</i>	8478	65	<i>now</i>	3081	115	<i>back</i>	1590
16	<i>what</i>	8171	66	<i>or</i>	3005	116	<i>time</i>	1580
17	<i>on</i>	7951	67	<i>did</i>	2911	117	<i>only</i>	1578
18	<i>is</i>	7816	68	<i>i don't</i>	2878	118	<i>you've</i>	1569
19	<i>have</i>	7802	69	<i>when</i>	2855	119	<i>off</i>	1555
20	<i>know</i>	7659	70	<i>had</i>	2829	120	<i>why</i>	1535
21	<i>one</i>	7488	71	<i>about</i>	2825	121	<i>something</i>	1510
22	<i>do</i>	7280	72	<i>want</i>	2823	122	<i>where</i>	1508
23	<i>was</i>	7133	73	<i>cos</i>	2796	123	<i>don't know</i>	1495
24	<i>got</i>	6842	74	<i>as</i>	2750	124	<i>could</i>	1486
25	<i>we</i>	6686	75	<i>mean</i>	2716	125	<i>she's</i>	1453
26	<i>he</i>	6618	76	<i>in the</i>	2662	126	<i>will</i>	1444
27	<i>don't</i>	6477	77	<i>my</i>	2504	127	<i>because</i>	1442
28	<i>they</i>	6475	78	<i>going</i>	2377	128	<i>have to</i>	1431
29	<i>but</i>	6178	79	<i>i mean</i>	2364	129	<i>you can</i>	1398
30	<i>so</i>	6148	80	<i>i've</i>	2327	130	<i>is it</i>	1390
31	<i>there</i>	6125	81	<i>put</i>	2303	131	<i>ah</i>	1380
32	<i>that's</i>	5957	82	<i>i think</i>	2286	132	<i>from</i>	1362
33	<i>for</i>	5673	83	<i>here</i>	2270	133	<i>his</i>	1358
34	<i>mm</i>	5662	84	<i>really</i>	2238	134	<i>if you</i>	1315
35	<i>not</i>	5270	85	<i>i'll</i>	2214	135	<i>nice</i>	1314
36	<i>go</i>	4941	86	<i>he's</i>	2212	136	<i>an</i>	1296
37	<i>be</i>	4869	87	<i>do you</i>	2196	137	<i>isn't</i>	1283
38	<i>this</i>	4781	88	<i>come</i>	2185	138	<i>mum</i>	1282
39	<i>get</i>	4772	89	<i>three</i>	2181	139	<i>what's</i>	1278
40	<i>like</i>	4744	90	<i>down</i>	2147	140	<i>thought</i>	1261
41	<i>just</i>	4696	91	<i>look</i>	2099	141	<i>any</i>	1254

42	<i>she</i>	4683	92	<i>didn't</i>	2074	142	<i>little</i>	1241
43	<i>all</i>	4459	93	<i>how</i>	2063	143	<i>of the</i>	1233
44	<i>er</i>	4441	94	<i>good</i>	2050	144	<i>and then</i>	1226
45	<i>yes</i>	4432	95	<i>you're</i>	2044	145	<i>more</i>	1220
46	<i>then</i>	4369	96	<i>there's</i>	2040	146	<i>haven't</i>	1212
47	<i>right</i>	4252	97	<i>gonna</i>	2007	147	<i>i don't know</i>	1204
48	<i>if</i>	4234	98	<i>her</i>	1983	148	<i>and i</i>	1192
49	<i>think</i>	4159	99	<i>some</i>	1950	149	<i>hundred</i>	1191
50	<i>can</i>	4148	100	<i>say</i>	1923	150	<i>much</i>	1180

There are 18 (12%) 2- or 3-grams among the top 150 items. This confirms the observation by O'Keeffe *et al.* (2006: 46) from their analysis of CANCODE concerning the high frequency of many chunks in spoken corpora. These items are marked in bold in Table 10. 12 of the 18 are in the third column of the list and thereby have a rank of 100 or greater. The first is *you know* at rank 58 with 3605 occurrences. The component words of this bigram are found at rank 2 (*you* 29688 occurrences) and rank 20 (*know* 7659 occurrences). The sole trigram in the top 150 items is *i don't know* with 1204 occurrences at rank 147. The unadjusted list should be compared with the adjusted frequency list in Table 11, where the procedure has been applied using a threshold value of five for both bigrams and trigrams for inclusion in the adjustment process. 43 (28.7%) of the top 150 items in the adjusted list are bi- or tri-gram items (marked in bold). The most frequent trigram in the adjusted list is *i don't know* (a move from rank 147 to rank 7). The three component words have experienced significant reduction: *i* (30371 [rank 1] → 660 [rank 17]), *don't* (6477 [rank 27] → 188 [rank 167]) and *know* (7659 [rank 20] → 51 [rank 1403]).

Table 11: Top 150 combined 1-, 2- and 3-grams in BNC Baby Demographic section after adjustment (using threshold for 2- and 3-grams of 5+ occs)

Rank	Item	Freq.	Rank	Item	Freq.	Rank	Item	Freq.
1	<i>yeah</i>	6877	51	<i>my</i>	398	101	<i>a bit of</i>	262
2	<i>mm</i>	3866	52	<i>a lot of</i>	395	102	<i>have a look</i>	258
3	<i>no</i>	3026	53	<i>that's</i>	392	103	<i>do you think</i>	255
4	<i>oh</i>	2002	54	<i>me</i>	380	104	<i>by</i>	254
5	<i>and</i>	1787	55	<i>isn't it</i>	373	105	<i>oh dear</i>	251
6	<i>yes</i>	1671	56	<i>this</i>	373	106	<i>on the</i>	249
7	<i>i don't know</i>	1204	57	<i>ha</i>	365	107	<i>four</i>	248
8	<i>the</i>	1103	58	<i>up</i>	364	108	<i>she's</i>	248
9	<i>what</i>	1100	59	<i>like</i>	361	109	<i>she</i>	246
10	<i>right</i>	887	60	<i>what do you</i>	356	110	<i>they</i>	246
11	<i>er</i>	808	61	<i>and the</i>	349	111	<i>down</i>	245
12	<i>erm</i>	800	62	<i>here</i>	338	112	<i>have you got</i>	244
13	<i>a</i>	754	63	<i>no no</i>	335	113	<i>are</i>	240
14	<i>in</i>	743	64	<i>please</i>	325	114	<i>three four five</i>	238
15	<i>that</i>	732	65	<i>who</i>	325	115	<i>first</i>	237
16	<i>well</i>	703	66	<i>i know</i>	324	116	<i>aye</i>	236
17	<i>i</i>	660	67	<i>just</i>	324	117	<i>good</i>	234
18	<i>or</i>	654	68	<i>that's right</i>	323	118	<i>aha</i>	232
19	<i>ah</i>	648	69	<i>anyway</i>	321	119	<i>him</i>	232
20	<i>it's</i>	630	70	<i>again</i>	319	120	<i>i think</i>	232
21	<i>then</i>	620	71	<i>out</i>	314	121	<i>where</i>	232
22	<i>it</i>	617	72	<i>today</i>	310	122	<i>for the</i>	231
23	<i>of</i>	591	73	<i>innit</i>	307	123	<i>of the</i>	228
24	<i>now</i>	568	74	<i>you have to</i>	307	124	<i>sorry</i>	228
25	<i>there</i>	564	75	<i>two</i>	306	125	<i>you've got to</i>	226
26	<i>you</i>	563	76	<i>eh</i>	304	126	<i>bloody</i>	224
27	<i>so</i>	546	77	<i>look</i>	304	127	<i>any</i>	223
28	<i>to</i>	544	78	<i>why</i>	304	128	<i>our</i>	223
29	<i>do you want</i>	543	79	<i>though</i>	300	129	<i>they're</i>	223
30	<i>for</i>	533	80	<i>yeah yeah</i>	297	130	<i>which</i>	223
31	<i>one two three</i>	520	81	<i>from</i>	296	131	<i>you know i</i>	221

32	mm mm	518	82	<i>he</i>	295	132	<i>pardon</i>	219
33	<i>i don't think</i>	517	83	<i>not</i>	295	133	<i>yep</i>	219
34	<i>ooh</i>	500	84	<i>as well</i>	294	134	<i>it's a</i>	217
35	<i>one</i>	491	85	<i>at</i>	294	135	<i>oh yes</i>	217
36	<i>is</i>	485	86	<i>them</i>	293	136	<i>their</i>	217
37	<i>mum</i>	482	87	<i>hello</i>	289	137	<i>you're</i>	217
38	<i>oh yeah</i>	481	88	<i>his</i>	286	138	<i>come on</i>	216
39	<i>really</i>	481	89	<i>was</i>	281	139	<i>daddy</i>	216
40	<i>but</i>	471	90	<i>you want to</i>	281	140	<i>probably</i>	216
41	<i>on</i>	470	91	<i>i'm</i>	278	141	<i>bye</i>	212
42	<i>two three four</i>	449	92	<i>her</i>	277	142	<i>some</i>	212
43	<i>is it</i>	447	93	<i>i mean i</i>	276	143	<i>these</i>	212
44	<i>with</i>	440	94	<i>as</i>	275	144	<i>with the</i>	212
45	<i>you know</i>	431	95	<i>off</i>	274	145	<i>oh no</i>	211
46	<i>your</i>	430	96	<i>dad</i>	273	146	<i>an</i>	210
47	<i>alright</i>	423	97	<i>actually</i>	272	147	<i>thank you</i>	207
48	<i>mhm</i>	411	98	<i>in the</i>	271	148	<i>we</i>	206
49	<i>mummy</i>	409	99	<i>he's</i>	267	149	<i>you know what</i>	206
50	<i>okay</i>	406	100	<i>no no no</i>	267	150	<i>that's it</i>	205

There is strong support, particularly in the case of *know*, for the claim that a standard (unadjusted) frequency list considerably inflates the frequency of single words that belong to larger chunks. Aside from *i don't know*, notice *you know* (rank 45), *i know* (rank 66), *you know I* (rank 131) and *you know what* (rank 149) as chunks containing *know*. In fact in the adjusted list there are 45 bi- and trigrams containing *know* with a higher rank than the single word item *know*. None of which, of course, were found above *know* in the unadjusted list.

Another interesting observation concerning differences between the unadjusted (Table 10) and adjusted (Table 11) frequency lists from BNC Baby Demographic is the rank reduction of many of the function words that routinely top any English frequency list. While the top ranking of personal pronouns *i* and *you* in the unadjusted list, above *the*, are an indication of spoken language, the top of the list is still quite generic. After adjustment, however, most of these items have dropped significantly in rank because of their participation in frequent chunks. The top of the adjusted list is now much more distinctly speech-like: *yeah*, *mm*, *oh*, *no*, *yes*, *right*. Further, many of the bi- and tri-gram chunks in the adjusted list are central clause fragments for questions (*do you want*, *what*

do you, do you think, have you got, do you know, can i have), directives (*have a look, you have to, you've got to*) or declarative statements (*i don't know, i don't think, i think it's, i want to*). Such observations need more detailed examination along with the application of the procedure to other spoken corpora.

6.2 BNC Baby Academic section

Table 12 shows the top 150 items from the frequency list for the academic section of the BNC Baby combining single words and 2- and 3-grams. 29 (19.3%) of the top 150 items are bi- or trigrams (marked in bold):

Table 12: Top 150 combined 1-, 2- and 3-grams according to type frequency in BNC Baby Academic section

Rank	Item	Freq.	Rank	Item	Freq.	Rank	Item	Freq.
1	the	70257	51	he	2219	101	there is	1071
2	of	44195	52	had	2134	102	social	1066
3	and	26672	53	that the	2132	103	work	1061
4	to	26245	54	than	2114	104	because	1053
5	in	25169	55	on the	2096	105	both	1044
6	a	23047	56	all	2075	106	is not	1042
7	is	17401	57	of a	2075	107	do	1035
8	that	12449	58	so	2019	108	may be	1027
9	of the	10454	59	his	2000	109	as the	1025
10	be	9364	60	i	1839	110	case	994
11	for	9303	61	its	1835	111	at the	984
12	as	9243	62	also	1784	112	now	983
13	it	8209	63	only	1765	113	history	980
14	are	7551	64	would	1759	114	3	978
15	by	7246	65	when	1730	115	used	970
16	this	7052	66	between	1728	116	even	969
17	with	6991	67	what	1709	117	being	964
18	in the	6244	68	for the	1668	118	has been	960
19	which	6026	69	no	1651	119	up	951
20	on	5779	70	by the	1642	120	have been	949
21	or	5508	71	you	1631	121	form	935

The adjusted frequency list: A method to produce cluster-sensitive frequency lists

22	not	5192	72	about	1597	122	very	932
23	was	5077	73	2	1507	123	each	921
24	have	4710	74	1	1501	124	species	921
25	from	4603	75	<i>can be</i>	1484	125	same	889
26	an	4491	76	where	1478	126	per	886
27	but	3890	77	who	1449	127	new	883
28	at	3612	78	<i>with the</i>	1447	128	must	882
29	we	3425	79	then	1431	129	<i>the same</i>	875
30	can	3238	80	two	1425	130	<i>in this</i>	872
31	<i>to the</i>	3225	81	any	1412	131	way	870
32	has	3159	82	<i>from the</i>	1407	132	law	862
33	there	3032	83	<i>in a</i>	1404	133	general	852
34	they	3001	84	into	1363	134	<i>this is</i>	841
35	were	2885	85	<i>as a</i>	1360	135	rather	838
36	their	2795	86	most	1354	136	<i>to a</i>	835
37	been	2790	87	those	1333	137	how	834
38	<i>it is</i>	2757	88	time	1255	138	b	831
39	more	2729	89	<i>is a</i>	1248	139	people	809
40	if	2705	90	them	1214	140	particular	808
41	one	2681	91	should	1205	141	much	807
42	formula	2635	92	example	1168	142	number	806
43	may	2619	93	however	1153	143	over	803
44	<i>and the</i>	2399	94	many	1146	144	might	800
45	<i>to be</i>	2384	95	<i>is the</i>	1127	145	within	800
46	such	2295	96	different	1122	146	data	797
47	these	2288	97	first	1115	147	see	791
48	will	2281	98	out	1101	148	given	782
49	other	2259	99	could	1093	149	does	757
50	some	2220	100	use	1074	150	thus	743

The unadjusted list should be compared with the list in Table 13 which is the adjusted frequency list using thresholds for 2- and 3-grams of 10+ and 20+ occurrences respectively. The use of two different thresholds here is simply to demonstrate how the parameters can be varied. Future work is required to ascer-

tain the most appropriate thresholds for different size corpora, different genres and different ranges of *n*. 40 of the top 150 items (26%) in the adjusted list are bi- or tri-grams compared to 29 of the top 150 in the unadjusted list (Table 12). This is a less marked change than the one seen with the Demographic section. Also there are no trigrams in the top 150 items either before or after adjustment.

Table 13: Top 150 combined 1-, 2- and 3-grams in BNC Baby Academic section after adjustment

Rank	Item	Freq.	Rank	Item	Freq.	Rank	Item	Freq.
1	<i>and</i>	12984	51	<i>when</i>	747	101	<i>other</i>	451
2	<i>the</i>	8500	52	<i>its</i>	723	102	<i>b</i>	446
3	<i>of</i>	7536	53	<i>can</i>	690	103	<i>like</i>	444
4	<i>in</i>	5523	54	<i>into</i>	689	104	<i>most</i>	444
5	<i>to</i>	4715	55	<i>l</i>	678	105	<i>for a</i>	442
6	<i>a</i>	4052	56	<i>if</i>	676	106	<i>while</i>	437
7	<i>or</i>	3903	57	<i>had</i>	672	107	<i>and a</i>	434
8	<i>for</i>	3414	58	<i>2</i>	668	108	<i>first</i>	429
9	<i>of the</i>	3225	59	<i>they</i>	661	109	<i>you</i>	429
10	<i>is</i>	3168	60	<i>we</i>	645	110	<i>here</i>	422
11	<i>by</i>	3094	61	<i>such</i>	635	111	<i>be</i>	421
12	<i>in the</i>	2705	62	<i>to a</i>	634	112	<i>without</i>	419
13	<i>are</i>	2667	63	<i>may</i>	627	113	<i>after</i>	417
14	<i>as</i>	2647	64	<i>then</i>	616	114	<i>work</i>	416
15	<i>with</i>	2330	65	<i>as the</i>	615	115	<i>thus</i>	413
16	<i>that</i>	2245	66	<i>can be</i>	605	116	<i>7</i>	412
17	<i>and the</i>	2110	67	<i>where</i>	605	117	<i>which is</i>	410
18	<i>was</i>	2044	68	<i>it is</i>	598	118	<i>but the</i>	406
19	<i>this</i>	1971	69	<i>only</i>	596	119	<i>often</i>	406
20	<i>which</i>	1948	70	<i>to be</i>	594	120	<i>if the</i>	404
21	<i>on</i>	1807	71	<i>however</i>	592	121	<i>c</i>	401
22	<i>from</i>	1719	72	<i>3</i>	588	122	<i>5</i>	400
23	<i>to the</i>	1680	73	<i>is the</i>	587	123	<i>about the</i>	400
24	<i>were</i>	1663	74	<i>also</i>	575	124	<i>10</i>	397
25	<i>an</i>	1445	75	<i>with a</i>	574	125	<i>species</i>	396

26	<i>but</i>	1405	76	is a	557	126	<i>her</i>	391
27	of a	1327	77	<i>who</i>	550	127	<i>very</i>	390
28	for the	1252	78	<i>any</i>	543	128	<i>data</i>	384
29	by the	1170	79	<i>both</i>	537	129	<i>6</i>	382
30	on the	1145	80	<i>so</i>	526	130	<i>many</i>	382
31	<i>their</i>	1113	81	<i>all</i>	525	131	<i>social</i>	380
32	from the	1100	82	<i>no</i>	520	132	should be	378
33	<i>it</i>	1087	83	<i>them</i>	520	133	have been	376
34	with the	1052	84	<i>4</i>	513	134	<i>through</i>	374
35	that the	1015	85	<i>formula</i>	509	135	<i>further</i>	373
36	<i>at</i>	985	86	<i>being</i>	505	136	<i>how</i>	371
37	in a	968	87	<i>would</i>	500	137	<i>see</i>	370
38	<i>these</i>	959	88	such as	497	138	<i>before</i>	369
39	<i>has</i>	953	89	<i>two</i>	496	139	for example	369
40	<i>his</i>	894	90	<i>not</i>	495	140	of an	368
41	<i>one</i>	881	91	<i>than</i>	495	141	will be	368
42	<i>some</i>	875	92	at the	487	142	<i>she</i>	367
43	<i>more</i>	863	93	<i>over</i>	484	143	<i>itself</i>	366
44	<i>have</i>	830	94	<i>now</i>	483	144	<i>should</i>	366
45	<i>about</i>	814	95	has been	472	145	<i>since</i>	363
46	<i>he</i>	805	96	may be	472	146	<i>could</i>	362
47	<i>will</i>	794	97	<i>p</i>	467	147	is not	361
48	<i>i</i>	792	98	<i>what</i>	465	148	<i>time</i>	361
49	as a	785	99	of this	455	149	of their	360
50	<i>between</i>	766	100	by a	452	150	<i>history</i>	358

Table 14 lists all of these chunks from the unadjusted and adjusted frequency lists with an indication of how their rank has changed after adjustment (↓ means reduction in rank [8 items], ↑ an increase [33 items] and = if it stays the same [1 item]). All of these items are bigrams so they have all dropped in frequency from the unadjusted list because the procedure begun with trigrams. The increase in rank of the majority of the items may be partially responsible for the drop in rank of the following items: *it is*, *to be*, *have been*, *is not*, *there is*, *the same*, *in this*, *this is*. However, each of these is part of a frequent trigram, e.g. *the*

same > the same time, the same as, the same way; is not > is not a, is not the, is not to, is not surprising; there is > there is a/an, there is no, there is little, there is some, there is also, there is evidence. In the case of *there is*, it is less frequent in the adjusted list (163 occurrences) than two 3-grams: *there is a* (301 occurrences) and *there is no* (245 occurrences).

Table 14: All bigrams in the top 150 items from 1-, 2- and 3-gram lists of BNC Baby Academic showing rank change after adjustment

item	change of rank?	rank before	rank after	frequency before	frequency after
<i>of the</i>	=	9	9	10454	3299
<i>in the</i>	↑	18	13	6244	2636
<i>and the</i>	↑	44	17	2399	2069
<i>to the</i>	↑	31	23	3225	1625
<i>for the</i>	↑	68	27	1668	1234
<i>of a</i>	↑	57	28	2075	1192
<i>by the</i>	↑	70	29	1642	1161
<i>from the</i>	↑	82	32	1407	1089
<i>on the</i>	↑	55	33	2096	1059
<i>that the</i>	↑	53	34	2132	1028
<i>with the</i>	↑	78	35	1447	1002
<i>in a</i>	↑	83	36	1404	998
<i>as a</i>	↑	85	49	1360	766
<i>to a</i>	↑	136	61	835	640
<i>it is</i>	↓	38	62	2757	636
<i>to be</i>	↓	45	67	2384	610
<i>as the</i>	↑	109	69	1025	596
<i>can be</i>	↑	75	70	1484	590
<i>with a</i>	↑	173	73	670	560
<i>is a</i>	↑	89	76	1248	540
<i>is the</i>	↑	95	77	1127	540
<i>such as</i>	↑	163	89	696	494
<i>has been</i>	↑	118	90	960	493

<i>may be</i>	↑	108	93	1027	480
<i>at the</i>	↑	111	98	984	455
<i>for a</i>	↑	209	102	568	439
<i>by a</i>	↑	241	103	500	433
<i>and a</i>	↑	280	105	444	431
<i>of this</i>	↑	155	109	719	429
<i>which is</i>	↑	205	115	572	408
<i>but the</i>	↑	309	119	406	403
<i>if the</i>	↑	231	122	520	399
<i>will be</i>	↑	177	126	661	384
<i>have been</i>	↓	120	130	949	378
<i>about the</i>	↑	297	131	418	377
<i>for example</i>	↑	158	132	713	375
<i>should be</i>	↑	227	141	531	366
<i>is not</i>	↓	106	147	1042	359
<i>there is</i>	↓	101	514	1071	163
<i>the same</i>	↓	129	166	875	333
<i>in this</i>	↓	130	178	872	318
<i>this is</i>	↓	134	177	841	319

In contrast to the adjusted list from the demographic texts (Table 11), the adjusted frequency list for academic writing has both a fewer number of chunks and lower ranked chunks with central clause functions. Instead we see more of the grammatical/discourse function chunks, such as *a number of*, *in terms of*, *as well as*, *for example*. In this case the procedure, at least using these values for N_{\max} and the frequency thresholds for n , has served to highlight the grammatical function items at the top of the list. The advantage is that it illustrates the frequent grammatical chunks (*of the*, *in the*, *with a*, *such as*, etc.) that one might tend to pass over when glancing over an n-gram list.

7 Final considerations and further developments

This paper is an initial attempt to address the recognized limitation in standard n-gram analysis when a range of values of n are combined, and particularly when single words and larger n-grams are combined into a single list. Because each size unit is counted on its own terms, the frequency for single words and

lower values of n will always be larger than (or perhaps equal to) the frequency of larger units. This was clearly illustrated in the analysis of spoken language from the BNC Baby corpus with the case of the word *know*.

The adjusted frequency list procedure is presented as one possible remedy for this problem. It gives priority to larger chunks (e.g. *on the other hand*) as it builds a frequency list by not counting the included components (*on the other, the other hand, on the, the other, other hand, on, the, other, hand*). Different results are achieved by varying the size of the largest n-gram (N_{\max}) at which the procedure begins and by applying different frequency (or potentially statistical) thresholds for the inclusion of specific n-grams. Three potential algorithms are presented here, two of which are used for implementation. There are likely to be other approaches as well. Future work is needed to apply the procedure to a range of corpora and to determine some criteria for determining appropriate thresholds.

The simple method presented here, along with other more complex techniques that have been recently proposed (Gries and Mukherjee 2010), demonstrates how corpus analysis continues to validate the importance of chunking in the investigation and description of language.

Notes

1. See Gries (2008) for a state-of-the-art overview of phraseological concepts in linguistic theory and computational method. The review addresses issues of terminological variation and the lack of generally accepted criteria for the identification and description of phraseologisms. Gries lays out six parameters/criteria that are designed to cover all of the aspects relevant to the notion of phraseology and by which any fully developed phraseological approach can be evaluated. He also makes suggestions with regards to methods and tools for the computational analysis of phraseology.
2. Sentence boundary punctuation has been observed in generating n-grams to produce the frequency list in Table 2. So for *...said the Mama bear. "Someone's been..."* the 3-gram *mama bear someone's* would not be counted. Not all software that produces n-gram lists respects sentence boundaries in this way.
3. Most often a number of iterations of the procedure will be required to capture the appropriate size for N_{\max} and the thresholds for a particular text or corpus. The procedure and suggested algorithms (see Section 4) are designed for such iterations. There are also more complex algorithms, such as the Lexical Gravity approach (most recently applied in Gries and

Mukherjee [2010]), that are designed to induce the maximum value of n for specific n -grams from the data.

References

- Ellis, Nick C. 1996. Sequencing in SLA. Phonological memory, chunking, and points of order. *Studies in Second Language Acquisition* 18: 91–126.
- Ellis, Nick C. 2003. Constructions, chunking, and connectionism: The emergence of second language structure. In C. Doughty and M.H. Long (eds.). *Handbook of second language acquisition*, 33–68. Oxford: Blackwell.
- Erman, Britt and Beatrice Warren. 2000. The idiom principle and the open choice principle. *Text* 20 (1): 29–62.
- Gries, Stefan Th. 2008. Phraseology and linguistic theory: A brief survey. In S. Granger and F. Meunier (eds.). *Phraseology: An interdisciplinary perspective*, 3–25. Amsterdam and Philadelphia: John Benjamins.
- Gries, Stefan Th. and Joybrato Mukherjee. 2010. Lexical gravity across varieties of English: An ICE-based study of n -grams in Asian Englishes. *International Journal of Corpus Linguistics* 15 (4): 520–548.
- Meunier, Fanny and Sylviane Granger (eds.). 2008. *Phraseology in foreign Language learning and teaching*. Amsterdam: John Benjamins.
- Nattinger, James R. and Jeanette S. DeCarrico. 1992. *Lexical phrases and language teaching*. Oxford: Oxford University Press.
- O'Donnell, Matthew Brook. 2008. KWICgrouper. Designing a tool for corpus-driven concordance analysis. In M. Scott, P. Pérez-Paredes and P. Sánchez-Hernández (eds.). *Software-aided analysis of language*. Special issue of *International Journal of English Studies* 8 (1): 107–122.
- O'Keeffe, Anne, Michael McCarthy and Ronald Carter. 2006. *From corpus to classroom: Language use and language teaching*. Cambridge: Cambridge University Press.
- Pawley, Andrew and Frances. H. Syder. 1983. Two puzzles for linguistic theory: Nativelike selection and nativelike fluency. In J. C. Richards and R. W. Schmidt (eds.). *Language and communication*, 191–226. London: Longman.
- Sinclair, John. 1991. *Corpus, concordance, collocation*. Oxford: Oxford University Press.
- Sinclair, John and Anna Mauranen. 2006. *Linear unit grammar: Integrating speech and writing*. Amsterdam: John Benjamins.

Stubbs, Michael. 2002. Two quantitative methods of studying phraseology in English. *International Journal of Corpus Linguistics* 7 (2): 215–244.

Wray, Alison. 2008. *Formulaic language: Pushing the boundaries*. Oxford: Oxford University Press.

Zimmer, Benjamin. 2010. On language: Chunking. *New York Times Sunday Magazine* September 19, 2010: 30.

Appendix

Text 1: Goldilocks and the Three Bears

(downloaded from http://www.dltk-teach.com/rhymes/goldilocks_story.htm)

Once upon a time, there was a little girl named Goldilocks. She went for a walk in the forest. Pretty soon, she came upon a house. She knocked and, when no one answered, she walked right in.

At the table in the kitchen, there were three bowls of porridge. Goldilocks was hungry. She tasted the porridge from the first bowl.

“This porridge is too hot!” she exclaimed.

So, she tasted the porridge from the second bowl.

“This porridge is too cold,” she said

So, she tasted the last bowl of porridge.

“Ahhh, this porridge is just right,” she said happily and she ate it all up.

After she’d eaten the three bears’ breakfasts she decided she was feeling a little tired. So, she walked into the living room where she saw three chairs. Goldilocks sat in the first chair to rest her feet.

“This chair is too big!” she exclaimed.

So she sat in the second chair.

“This chair is too big, too!” she whined.

So she tried the last and smallest chair.

“Ahhh, this chair is just right,” she sighed. But just as she settled down into the chair to rest, it broke into pieces!

Goldilocks was very tired by this time, so she went upstairs to the bedroom. She lay down in the first bed, but it was too hard. Then she lay in the second bed, but

it was too soft. Then she lay down in the third bed and it was just right. Goldilocks fell asleep.

As she was sleeping, the three bears came home.

“Someone’s been eating my porridge,” growled the Papa bear.

“Someone’s been eating my porridge,” said the Mama bear.

“Someone’s been eating my porridge and they ate it all up!” cried the Baby bear.

“Someone’s been sitting in my chair,” growled the Papa bear.

“Someone’s been sitting in my chair,” said the Mama bear.

“Someone’s been sitting in my chair and they’ve broken it all to pieces,” cried the Baby bear.

They decided to look around some more and when they got upstairs to the bedroom, Papa bear growled, “Someone’s been sleeping in my bed,”

“Someone’s been sleeping in my bed, too” said the Mama bear

“Someone’s been sleeping in my bed and she’s still there!” exclaimed Baby bear.

Just then, Goldilocks woke up and saw the three bears. She screamed, “Help!” And she jumped up and ran out of the room. Goldilocks ran down the stairs, opened the door, and ran away into the forest. And she never returned to the home of the three bears.

