

Harnessing keyness: Corpus-based approach to ESP material development



John Blake

Background

ESP course developers adopting a corpus-based or corpus-driven approach¹ can create a **focus corpus** of relevant texts and use **keyness**² to identify **lexical sets** to integrate into teaching materials.

Keyness

- Keyness is a measure of the frequency of disproportionate occurrence.
- Focus corpus is contrasted to a reference corpus.
- Operational definition of keyness³ affects key words.
- Words may be key to whole corpus (global) or part (localised or bursty)⁴.

This study investigates the effect of different **1**. concordancers, **2**. reference corpora and 3. statistical formulae on the keyword lists for a focus corpus of research articles on international business.



Table 1: Focus corpus

Variable	Count
Tokens	2,516,051
Words	1 966 650

	±,500,05
Sentences	77,547

77,547

1: Concordancer

Concordancers are designed for different purposes and budgets. Raw frequency counts also differ with concordancer⁵. 4th generation concordancers⁶ are webbased, fast and suitable for large corpora. This presentation compares the results using two popular concordancers: AntConc and Sketch Engine.

2.3w (Windows) 2011 Concordance Plot File View Clusters Collocates Word List Keyword List 1.sts 2.sts 1.sts 2.sts 1 the advance of bospice is allowing more and more terminally taa_unit2.sts 1 the advance of bospice is allowing more and more terminally taa_unit2.sts 1 the advance of bospice is allowing more and more terminally taa_unit2.sts 1 the advance of bospice is allowing more and more terminally taa_unit2.sts 2 the advance of bospice is allowing more and more terminally taa_unit2.sts 2 the advance of bospice is allowing more and more terminally taa_unit2.sts 3 the advance of bospice is allowing more and more terminally taa_unit2.sts 3 the advance of bospice is allowing more and more terminally taa_unit2.sts 3 the second provide the second provide terminal term	Table 2: Examples of 3G &4G concordancers		Raw frequency list Sketch Engine 1 the 106,022		Raw frequency list AntConc 1 the 106,064		
Print programs. And it is expanding the materials it will der the standist to the second s	3G	AntConc ⁷ UAM Corpus Tool ⁸ Wordsmith Tools ⁹	 2 and 3 of 4 to 5 in 6 a 	77,508 72,733 47,454 41,791 32,007	2 and 3 of 4 to 5 in 6 a	x 77,542 72,990 47,834 42,056 32,336	
Fig 1: Sceenshot of AntConc	4G	CQPweb ¹⁰ Sketch Engine ¹¹ Wmatrix ¹²	7 that8 is9 for10 as	23,092 21,249 17,293 14,309	7 tha8 is9 for10 as	t 23,092 21,245 17,303 14,329	
			L	ist 1		List 2	

2: Reference corpus

The **genre** & **diachrony** of a corpus significantly affect keyness^{13, 14}.

Table 3: Comparison of Brown and BAWE			
Brown cornus	British Acadomic Writton		

Keyword list Sketch Engine Brown Midway (1000) 1 firms

Keyword list Sketch Engine BAWE **Midway (1000)** firms

AntConc users need to provide their own reference corpus, while Sketch Engine provides access to 20 corpora, including enTenTen12, a 12 billion token corpus.

DIOWII COLPUS	DITUSTI ACQUEITIC VVITUETI	2 firm	2 firm	
	English Corpus (BAWE)	3 export	3 export	
circa 1960s	circa 2000s	4 foreign 5 subsidiary	4 table 5 variables	
General corpus	Academic corpus	6 internationalization 7 FDI	6 international 7 markets	
American English	British English	8 subsidiaries	8 knowledge	
1.000.000 words	6.506.995 words	9 markets 10 MNEs	9 foreign 10 market	
		List 3	List 4	

3: Statistical formula

Many linguists are ill at ease with statistics¹⁵.

Main ideas

- Frequency bias of log likelihood and chi squared³
- Rarity bias of **mutual information**
- The simple maths version¹⁶ in Sketch • Engine names the formulae clearly and assumes language is not random¹⁷.

Keyword list	Keyword list	Keyword list	Keyword list	Keyword list
AntConc	AntConc	Sketch Engine	Sketch Engine	Sketch Engine
Brown	Brown	BAWE	BAWE	BAWE
Log likelihood	Chi squared	Rare (0.01)	Midway (1000)	Common (1 million)
1 the	1 the	1 OFDI	1 firms	1 and
2 firms	2 firms	2 offshoring	2 firm	2 firms
3 firm	3 firm	3 Vahlne	3 export	3 firm
4 al	4 et	4 multinationality	4 table	4 foreign
5 et	5 al	5 full-size	5 variables	5 knowledge
6 in	6 in	6 MathML	6 international	6 international
7 knowledge	7 knowledge	7 Kogut	7 markets	7 market
8 market	8 market	8 BOP	8 knowledge	8 country
9 this	9 international	9 MathJAx	9 foreign	9 table
10 table	10 foreign	10 Ghoshal	10 market	10 performance
List 5	List 6	List 7	List 8	List 9

Summary

4G concordancers with the choice of reference corpora and formulae enable developers to tweak results to create the most useful list of keywords.

References

1. Tognini-Bonelli, E. (2001). Corpus Linguistics at Work. Amsterdam: John Benjamins. 2. Scott, M. (1997). PC analysis of key words – and key key words. System, 25 (1), 1-13. 3. Gabrielatos, C. and Marchi, A. (2012) Keyness: Appropriate metrics and practical issues. Paper presented at Corpus-assisted Discourse Studies International Conference 2012. University of Bologna, Italy. 13-14 September, 2012. 4. Katz, S. (1996). Distribution of Common Words and Phrases in Text and Language

Modelling, Natural Language Engineering, 2 (1), 15-59.

5. Anthony, L. (2012). A critical look at software tools in corpus linguistics. Linguistic *Research, 30* (2), 141-161.

6. McEnery, T, & Hardie, A. (2012). Corpus linguistics: Method, Theory and Practice. Cambridge: Cambridge University Press.

7. Anthony, L. (2012). AntConc (Version 3.2.4) [Computer Software]. Tokyo, Japan: Waseda University.

8. O'Donnell, M. (2013). UAM Corpus Tool (Versions 2.8 & 3.1). Wagsoft Systems. 9. Scott, M. (2012). WordSmith Tools (Version 6). Liverpool: Lexical Analysis Software. 10. Hardie, A. (2012). CQPweb - combining power, flexibility and usability in a corpus analysis tool. International Journal of Corpus Linguistics 17 (3), 380–409.

11. Kilgarriff, A. et al. (2004). The Sketch Engine. Lexical Computing.

12. Rayson, P. (2008). W-matrix corpus analysis and comparison tool. Lancaster University. 13. Scott, M. (2009). In search of a bad reference corpus. In D. Archer (ed.), What's in Word*list? Investigating Word Frequency and Keyword Extraction (pp.79-92).* Oxford: Ashgate. 14. Goh, G-Y. (2010). Choosing a reference corpus for keyword extraction. *Linguistic Research*, 28 (1), 239-256.

15. Loewen, S. et al. (2014). Statistical literacy among applied linguists and second language acquisition researchers. TESOL Quarterly, 48 (2). 360-388.

16. Kilgarriff, A. (2009). Simple maths for keywords. In Mahlberg, M., González-Díaz, V. & Smith, C. (eds.), Proceedings of the Corpus Linguistics Conference CL2009. University of Liverpool, UK, 20-23 July 2009.

17. Kilgarriff, A. (2005). Language is never ever ever random. Corpus Linguistics and Linguistic *Theory* 1 (2): 263-276.