

This was an invited chapter for a book dealing with the history of English, which was published by a major university press in 2015. After two positive reviews, it was removed at the last minute because of objections from a third reviewer. This person apparently had a vested interest in preventing others from seeing [Section 2](#), which s/he felt was too explicit in showing the limitations of small corpora for lexical studies.

I believe that this paper explains some very important concepts about corpora and lexical change, which are often misunderstood (or willfully ignored and even suppressed) by other researchers. As a result, I've decided to make the paper freely available via the Web – where it will likely have wider circulation than it would have had, sitting on the shelf in a library.

## **Using corpora to research lexical changes in Late (and Early) Modern English**

Mark Davies  
Brigham Young University

The focus of this chapter is the different ways in which corpora English can be used to research lexical changes in Late (and Early) Modern English. Although we could examine any type of language change (syntactic, morphological, pragmatic, and so on), in this case we have focused on lexical change, because it is an often-overlooked area of change in English – especially by those who are used to working with small corpora, where there are typically not enough tokens to look at medium and low-frequency words. In addition, we will focus primarily on Late Modern English – especially the period from the early 1800s through the mid-1900s – because until recently there were very few large, structured corpora for this period. As we will see in this chapter, however, there are a number of corpora and corpus-related resources that have become available in just the past four or five years, which have greatly expanded our ability to carry out extremely detailed investigations on lexical changes in English.

We will begin this chapter by looking at the type of research that one can do with the 400 million word Corpus of Historical American English (COHA). We will then consider what type of lexical research we can (or cannot) do with smaller corpora, a handful of large text archives, the quotations in the online Oxford English Dictionary (OED), and Google Books.

As we consider all of these resources, we will focus on how the corpora and related resources can be used to generate the following types of data:

1. Concordances for a particular word or phrase, to see the patterns in which it occurs
2. The frequency of specific words and phrases over time
3. Finding all words that are more frequent in one period than in another
4. Collocates of a particular word or phrase (over time) in order to see semantic change
5. Easy access to the relative frequency of related words, such as synonyms or words in a user-defined list

# 1. The Corpus of Historical American English (COHA)

The Corpus of Historical American English (COHA) was released in 2010 as a “companion” corpus to the 450 million word Corpus of Contemporary American English (COCA), which contains texts from just 1990 to the current time. COHA, on the other hand, contains 400 million words from 1810-2009, which makes it at least 100 times as large as any other publicly-available, structured corpus of English (as opposed to just text archives). It is composed of texts from fiction, popular magazines, newspapers, and non-fiction books, and it is genre-balanced, in the sense that it has roughly the same balance of these genres from decade to decade. (For an overview of COHA, see Davies 2012a and 2012b).

Most importantly – for the purpose of this paper – it allows for an extremely wide range of research on lexical changes in Late Modern English, which in many cases are not available from other corpora. Due to limited space, we will provide just one sample of each type of search, but extensive tutorials, with many additional sample searches, can be found at the corpus website.

At the most basic level, we can use COHA to see concordances for a particular word or phrase, to see the patterns in which it occurs. Users simply enter a word, phrase, or construction, and then (optionally) select a range of years, to see up to 1000 entries at a time. For example, at the current time, *fathom* as a verb is used almost exclusively in negative and interrogative contexts, and COHA shows that the same holds true when we examine the concordance lines from the 1820s-1880s, e.g.:

Figure 1. COHA: concordance lines for *fathom*, 1820s-1880s

32	1845	FIC	FleetwoodTheStain	A B C	a mystery in Fleetwood's interrogation, which <b>she could not fathom</b> . The moment was a crisis in her fate. A straw
33	1857	FIC	HabeVaughan	A B C	watchful eyes of her father and aunt; <b>she dared not fathom</b> her own unhappiness; but, continuing her customary round of
34	1844	FIC	Bondmaid	A B C	counsels of the eternal powers, and <b>she could not fathom</b> them: <b>What</b> the breast of the Almighty, lies many a
35	1835	FIC	AdventuresTimothy	A B C	where is the human philosophy that can <b>explain the operation of</b> fathom its causes? If I rightly understand the history of these cases
36	1859	FIC	TrueWomanhoodA	A B C	which I do not understand – which I <b>can not possibly fathom</b> . He knows too much about us – and much more,
37	1882	FIC	Poems	A B C	I combined -- Ay! shrewd scientist <b>too</b> <b>she shall fathom</b> your <b>mind</b> , <b>shall</b> plumb that strange sea to the uttermost deep

Concordance lines can also provide evidence for semantic change over time. For example, the following two figures are for the word *gay* in the 1860s-1910s and the 1970s-2000s:

Figure 2. COHA: concordance lines for *gay*, 1860s-1910s

30	1905	FIC	RainbowRose	A B C	eyes see Such buds of scented joys to <b>be</b> A <b>gay</b> <b>green</b> <b>garden</b> <b>so</b> softly fanned By the blithe breeze that blows To
31	1885	FIC	HomeScenesHome	A B C	wife riding out in a handsome vehicle, <b>drawn by a</b> <b>gay</b> <b>horse</b> ; <b>and</b> <b>taking</b> their comfort, " said I, as
32	1896	FIC	HolidayStories	A B C	Aunt Hetty, in a blue gingham <b>down</b> , with a <b>gay</b> <b>kerchief</b> <b>red</b> <b>on</b> her head, was slowly and pensively rocking
33	1911	FIC	GraceHarloweJunior	A B C	with light. There was a hum of <b>gay</b> <b>voices</b> and <b>gay</b> <b>laughter</b> <b>and</b> all the pleasant excitement attending an amateur
34	1902	FIC	Audrey	A B C	and Nicholson and Duke of Gloucester <b>streets were blown</b> the <b>gay</b> <b>leaves</b> <b>of</b> <b>the</b> mornings white frosts lay upon the earth like
35	1885	MAG	Atlantic	A B C	derkey, and lumbering coach, and <b>throng of</b> <b>gay</b> <b>life</b> <b>have</b> <b>disappeared</b> . There was no room in this valley

Figure 3. COHA: concordance lines for *gay*, 1970s-2000s

20	2002	MAG	WashNorth	A B C	, but it continues to puzzle me: <b>Clearly</b> , the <b>gay</b> <b>community</b> <b>is</b> <b>a</b> <b>whole</b> is supportive of safe sex, since the
21	1999	NF	DeadlyPersuasion	A B C	general population.) // However, many <b>advertisers who target</b> <b>gay</b> <b>consumers</b> <b>are</b> <b>likely</b> to remain closeted about it, for fear of
22	2004	MAG	Time	A B C	next month, to come up with a <b>solution that protects</b> <b>gay</b> <b>couple's</b> <b>rights</b> . (Bawden suggested civil unions.) If
23	1983	FIC	PlayExecutionJustice	A B C	of San Francisco, George Moscone, and the first <b>openly</b> <b>gay</b> <b>election</b> <b>official</b> City Supervisor Harvey Milk, reads like a Who
24	2003	FIC	AllIWantIsEverything	A B C	heart-wrenching sob story about how her <b>mother had divorced</b> her <b>gay</b> <b>father</b> <b>and</b> <b>was</b> <b>about</b> to marry a man she barely knew,
25	1993	NEWS	Atlanta	A B C	, it is mind-boggling to imagine <b>what activities constitute</b> <b>gay</b> <b>games</b> . # County's gays feel isolated # Commission

Second, COHA allows users to quickly and easily see the frequency of specific words and phrases over time. Figures 1-3 below are just a handful of examples, showing the frequency of

*teenager* (increasing over time), *steamship* (highest frequency in the early 1990s, and related to changes in society and culture), and *many a time* (decreasing over time).

Figure 4. COHA: *teenager*

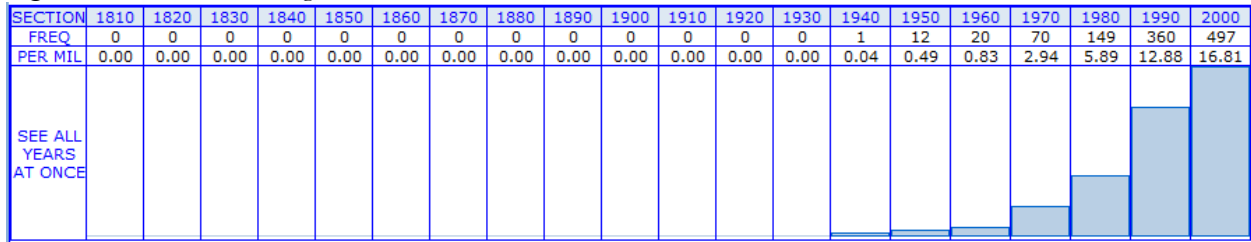


Figure 5. COHA: *steamship*

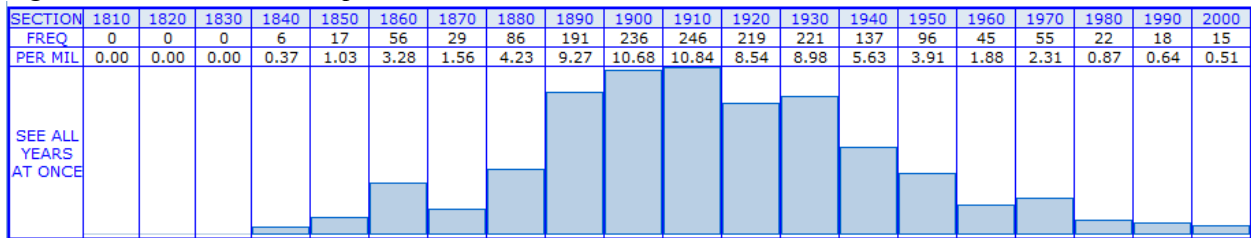
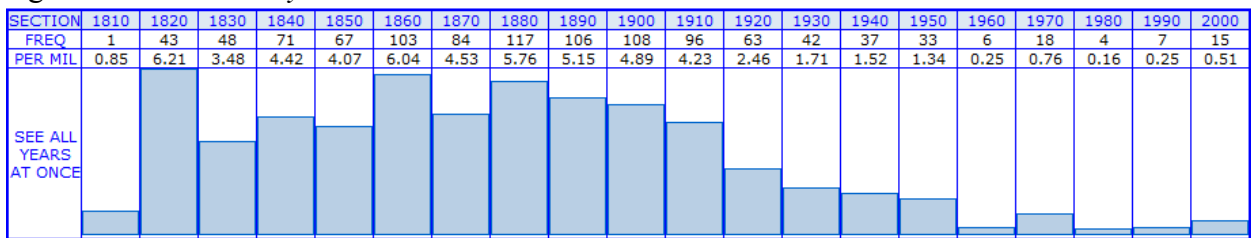
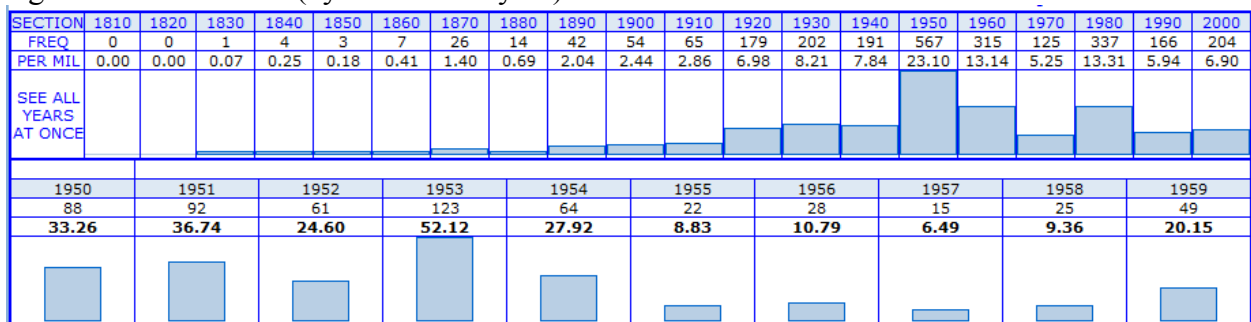


Figure 6. COHA: *many a time*



It is also possible to see the frequency by individual year, as in the following chart for *reds*, which peaks in 1953 (the height of the McCarthy anti-Communist hearings in the US Senate):

Figure 7. COHA: *reds* (by decade and year)



Users can also see the frequency (by decade) of each matching string for a particular construction. For example, the following table shows the first few entries for the construction *many a* [noun] (for an in-depth discussion of *many a* NOUN in COHA, see Hilpert 2012):

Figure 8: COHA: *many a* [noun]

CONTEXT		ALL	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000
1	MANY A TIME	1069	1	43	48	71	67	103	84	117	106	100	96	83	42	37	33	6	18	4	7	15
2	MANY A DAY	530		6	34	42	40	47	39	64	64	48	30	36	19	15	9	3	6	4	3	1
3	MANY A MAN	455	1	4	21	25	32	22	32	24	47	58	60	42	29	13	11	12	10	10	3	3
4	MANY A YEAR	435		17	18	36	39	39	29	42	55	38	21	30	15	16	12	0	5	6	5	5
5	MANY A NIGHT	181	1	3	8	9	10	18	14	12	15	11	15	14	14	9	6	1	7	3	6	3
6	MANY A WOMAN	104		5	1	4	5	10	9	8	11	13	9	12	5	4	3		2	1	1	1

Users can then click on the bar for any decade or year in the frequency chart display (Figures 4-7) or words and phrases in the table display (Figure 8) to see the KWIC entries (Keyword in Context); for example the entries for *many a time* in the 1850s:

Figure 9. COHA: Keyword in Context lines (KWIC)

26	1853	FIC	UncleSamsEmancipation	A	B	C	Druids -- graybeards, dusky garments and all, on the shores of Maine, <b>many a time</b> ; and if anybody wants to feel the beauty and grandeur of the
27	1854	FIC	TempestSunshine	A	B	C	, and which Julia, in the bitterness of her heart, cursed many and <b>many a time</b> . In the early part of the morning Dr. Lacey wandered down to
28	1854	FIC	TempestSunshine	A	B	C	had resolved that if there still was hope for him he would find it. <b>Many a time</b> during the succeeding days he prayed in secret, not that Fanny might
29	1854	FIC	TempestSunshine	A	B	C	the bookcase the old family Bible, on whose dark dusty covers she remembered having <b>many a time</b> written her name. All was now explained. Her father's
30	1854	FIC	RhymesWithReason	A	B	C	storm. " Behold upon the kitchen wall That old and rusty gun! Full <b>many a time</b> the same I've borne T'ill setting of the sun; On training-days
31	1854	FIC	ThisThatOther	A	B	C	touches to a beautiful Madonna. Wearily had the girl-artist toiled and studied, and <b>many a time</b> had her lamp grown dim, in the gray light of morning,
32	1854	FIC	MartinHerrysale	A	B	C	's disconsolate soul, his step-mother was remembered with a sort of tenderness, and <b>many a time</b> he wished himself sitting in her kitchen corner, as of old
33	1855	FIC	ScenesCharacters	A	B	C	Reginald's. It was not very often that quarrels went so far, but <b>many a time</b> in thought, word, and deed was the rule of love transgressed

Third, with one simple search, in COHA it is possible to see all words that are more frequent in one period than in another. For example, the following table compares *-ism* words in the 1860s-1910s (left) with the 1970s-2000s (right), and it is quite interesting to see how these relate to cultural changes in the United State. Note that users could also compare any set of words – for example, all verbs or all nouns that are more common in one period than another.<sup>1</sup>

Figure 10. COHA: *-ism* words, 1860s-1910s vs 1970s-2000s

SEC 1: 121,332,176 WORDS						SEC 2: 106,640,094 WORDS					
WORD/PHRASE	TOKENS 1	TOKENS 2	PH 1	PH 2	RATIO	WORD/PHRASE	TOKENS 2	TOKENS 1	PH 2	PH 1	RATIO
1 PAUPERISM	217	1	1.79	0.01	190.72	1 RACISM	994	0	9.32	0.00	932.11
2 FETICHISM	113	0	0.93	0.00	93.33	2 TOURISM	759	0	7.12	0.00	711.74
3 BIMETALLISM	62	0	0.51	0.00	51.10	3 ACTIVISM	405	1	3.80	0.01	460.85
4 ROMANTISM	62	0	0.51	0.00	51.10	4 MARXISM	342	2	3.21	0.02	194.56
5 HEATHENISM	94	2	0.77	0.02	41.31	5 FUNDAMENTALISM	176	0	1.65	0.00	165.04
6 DEMAGOGISM	41	1	0.34	0.01	36.04	6 MULTICULTURALISM	156	0	1.48	0.00	146.29
7 PROPAGANDISM	43	0	0.35	0.00	35.44	7 COUNTERTERRORISM	134	0	1.26	0.00	125.66
8 ECCLESIASTICISM	42	0	0.35	0.00	34.82	8 DYNAMISM	130	0	1.22	0.00	121.91
9 MOHAMMEDANISM	117	3	0.96	0.03	34.28	9 AUTHORITARIANISM	101	1	0.95	0.01	114.92
10 MONOPOLISM	40	0	0.33	0.00	32.97	10 EXPRESSIONISM	100	1	0.94	0.01	113.78
11 SPIRITISM	38	0	0.31	0.00	31.32	11 CONSUMERISM	121	0	1.13	0.00	113.47
12 INVALIDISM	70	2	0.58	0.02	30.76	12 SEXISM	116	0	1.09	0.00	108.78

Fourth, in COHA we can see the collocates of a particular word or phrase (over time) in order to see semantic change. For example, to continue the example shown above, the following chart shows the collocates of *gay* by decade. The older meaning of “happy, cheerful” is seen in lines like 1 and 2, while the newer meaning related to sexual orientation is found in lines like 4 and 6.

<sup>1</sup> To see the actual query online, see <http://corpus.byu.edu/coca/?c=coca&q=20252401>. These are adjectives that occur at least two times in each of the two time periods. The [ratio] column shows how much more common the words is in the one time period than the other, taking into account the overall size (in words) of each time period.



Figure 11. COHA: collocates of *gay*, by decade

	CONTEXT	ALL	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000	
1	BRIGHT	172	1	5	0	10	14	13	23	12	14	12	4	12	12	8	11	7	4	2			
2	HAPPY	153		2	13	14	7	19	8	9	11	8	12	11	14	3	8	8	3	1	2		
3	FLOWERS	152		5	13	10	17	9	18	16	7	13	10	11	7	5	6	1	3			1	
4	RIGHTS	128																	6	19	47	57	
5	COLORS	127		3	6	4	9	13	8	9	10	5	7	10	6	17	8	3	6	1			
6	LESBIAN	117										1							1	3	49	63	
7	LAUGH	112		2	4	4	14	12	8	14	4	8	12	8	10	2	4	4	4				
8	MARRIAGE	93				1		1	1					1					1			7	81

It is also possible to compare the collocates in different periods. For example, the following table shows the collocates of *gay* in the 1830s-1910s (left) compared to the 1970s-2000s (right):

Figure 12. COHA: ADJ/NOUN collocates near the noun *gay*

SEC 1: 167,626,806 WORDS						SEC 2: 106,640,894 WORDS							
	WORD/PHRASE	TOKENS 2	TOKENS 1	PH 2	PH 1	RATIO		WORD/PHRASE	TOKENS 2	TOKENS 1	PH 2	PH 1	RATIO
1	GRAVE	63	1	0.38	0.01	40.08	1	LESBIAN	116	1	1.05	0.01	182.34
2	LADY	66	0	0.39	0.00	39.37	2	RIGHTS	129	0	1.21	0.00	120.97
3	GALLANT	65	0	0.39	0.00	38.79	3	COMMUNITY	81	2	0.76	0.01	63.66
4	BRIGHT	57	0	0.34	0.00	34.00	4	BAR	38	1	0.36	0.01	59.73
5	HEART	51	1	0.30	0.01	32.44	5	STRAIGHT	38	1	0.33	0.01	35.02
6	GRAVE	51	0	0.30	0.00	30.42	6	LESBIAN	35	1	0.33	0.01	55.02
7	THROW	50	0	0.30	0.00	29.83	7	ACTIVISTS	31	1	0.29	0.01	48.73
8	VOICES	46	1	0.27	0.01	29.26	8	MARRIAGE	89	3	0.83	0.02	46.63
9	GLAD	49	0	0.29	0.00	29.23	9	COUPLES	27	1	0.25	0.01	42.44
10	SPIRITS	48	0	0.29	0.00	28.64	10	NATIONAL	20	1	0.19	0.01	31.44
11	SONG	46	0	0.27	0.00	27.44	11	BISEXUAL	31	0	0.29	0.00	29.07
12	ATTIRE	42	1	0.25	0.01	26.72	12	LESBIANS	27	0	0.25	0.00	25.32

Another example of comparing collocates is Figure 13, which shows the adjectival collocates preceding *women* in the 1830s-1890s (left) and the 1960s-2000s (right), and which shows how women are represented and portrayed in the two periods:

Figure 13. COHA: Adjectival collocates of *women*

SEC 1: 122,828,575 WORDS						SEC 2: 130,617,326 WORDS							
	WORD/PHRASE	TOKENS 2	TOKENS 1	PH 2	PH 1	RATIO		WORD/PHRASE	TOKENS 2	TOKENS 1	PH 2	PH 1	RATIO
1	STRONG-MINDED	23	1	0.19	0.01	24.46	1	PREGNANT	264	3	2.02	0.02	82.79
2	NOBLE	30	2	0.24	0.02	15.95	2	BATTERED	70	0	0.54	0.00	53.99
3	TRUE	13	1	0.11	0.01	13.82	3	NATIONAL	70	0	0.54	0.00	53.99
4	AGED	24	2	0.20	0.02	12.76	4	AFRICAN-AMERICAN	61	0	0.47	0.00	46.70
5	CULTIVATED	12	0	0.10	0.00	9.77	5	PROFESSIONAL	47	1	0.36	0.01	44.20
6	DEFENCELESS	11	0	0.09	0.00	8.96	6	JAPANESE	39	1	0.30	0.01	36.57
7	ELDER	11	0	0.09	0.00	8.96	7	BLACK	493	14	3.77	0.11	33.11
8	FINEST	8	1	0.07	0.01	8.31	8	NAKED	69	2	0.53	0.02	32.44
9	LOVELIEST	8	1	0.07	0.01	8.51	9	SOVIET	32	0	0.24	0.00	24.50
10	LITERARY	10	0	0.08	0.00	8.14	10	CATHOLIC	25	1	0.19	0.01	23.51
11	HEROIC	10	0	0.08	0.00	8.14	11	DIVORCED	28	0	0.21	0.00	21.44
12	DEVOTED	10	0	0.08	0.00	8.14	12	HOMELESS	22	1	0.17	0.01	20.69

Fifth and finally, in COHA it is possible to see the relative frequency of related words, such as synonyms or words in a user-defined list. For example, the following table shows the frequency of synonyms of *dangerous* by decade, where we see the increase of words like *critical* and *risky*, and the decrease of words like *perilous*, *alarming*, and *treacherous*:

Figure 14. COHA: Frequency of synonyms of *dangerous* by decade

CONTEXT	ALL	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000
1 SERIOUS [S]	42404	78	633	908	1094	1331	1311	1714	2141	2214	2346	2452	2661	2544	2566	2544	2547	2126	2248	3186	3197
2 DANGEROUS [S]	24536	107	519	892	942	986	1074	1197	1315	1287	1225	1308	1417	1355	1447	1442	1442	1212	1696	1667	1811
3 CRITICAL [S]	14238	51	171	313	307	350	400	472	485	576	708	768	705	618	903	880	1119	1124	1241	1350	1603
4 GRAVE [S]	10638	18	214	432	552	627	710	763	854	748	763	637	725	676	627	552	687	361	388	210	241
5 DARING [S]	4722	46	182	360	369	262	276	241	289	294	423	342	302	312	161	170	188	185	108	124	154
6 PERILOUS [S]	3121	11	157	270	254	259	252	230	199	179	192	170	192	114	105	96	85	93	92	89	76
7 ALARMING [S]	3080	10	142	184	156	129	181	158	181	117	139	125	180	141	160	133	130	138	158	190	162
8 TREACHEROUS [S]	2776	27	121	183	191	129	176	192	191	158	134	168	158	123	101	114	99	82	114	125	152
9 HAZARDOUS [S]	2221	6	85	108	81	98	72	75	85	82	72	81	118	87	104	95	115	117	150	426	185
10 PRECARIOUS [S]	1871	18	63	70	76	70	81	81	66	82	83	123	148	123	113	130	119	112	107	110	103
11 RISKY [S]	1350		4		1	2	10	26	34	35	46	66	49	58	87	88	105	137	180	243	372

Another example is the search [=beautiful] woman, which finds any synonym of beautiful followed by woman (note the decrease in *lovely* and *charming* and the increase in *attractive*, *striking*, and *gorgeous*):

Figure 15. COHA: Synonyms: “beautiful” + woman

CONTEXT	ALL	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000
1 BEAUTIFUL WOMAN	1321		18	30	97	10	38	88	71	76	94	48	71	52	33	70	63	74	40	106	162
2 HANDSOME WOMAN	214		2	2	8	5	10	39	30	20	24	13	18	14	20	47	18	12	21	8	18
3 LOVELY WOMAN	218	1	6	11	40	20	18	33	30	24	16	18	7	15	12	14	7	8	16	9	9
4 ATTRACTIVE WOMAN	226			1	2		2	5	4	6	7	3	9	11	20	11	17	26	32	22	48
5 WONDERFUL WOMAN	164		3	4	2	1	3	12	8	7	7	21	17	11	16	5	7	12	13	8	9
6 CHARMING WOMAN	145		3	2	3	6	8	19	14	13	17	9	19	11	4	5	2	4			2
7 GOOD-LOOKING WOMAN	84			1	3	1	3	3	2	1	1	1	7	9	4	1	3	12	11	13	13
8 MAGNIFICENT WOMAN	34			2	7	2	1	1	4		2		4	1	2		1		3	1	2
9 FINE-LOOKING WOMAN	25			1		2	1	1	4	1	3	3	1	1	1		2		3	1	
10 STRIKING WOMAN	19										2				1	1	2		4	5	4
11 GORGEOUS WOMAN	18															3		4	2	6	

We have shown how COHA – which extends from the 1810s-2000s, can be used to examine a wide range of lexical changes in English. As is discussed in Davies (2011), it is possible to use the related Corpus of Contemporary American English (currently 450 million words, 1990-2012) to carry out similar searches for just the last two decades or so. For example, we could see charts (like Figures 1-3 above) showing the frequency by year of words and phrases that have increased in frequency during this time (e.g. *old-school*, *morph* (verb), *freak out*, *perfect storm*, *think outside the box*, or *throw someone under the bus*).

As we can do for COHA for the period of the 1810s-2000s, we can use COCA to look at just the 1990s-2000s to generate lists of all words that have increased or decreased in frequency during this time. For example, with one quick search we could quickly find all adjectives that have decreased in frequency from the 1990s to the 2000s (e.g. *Croatian*, *Kuwaiti*, *Bosnian*, *US-Soviet*, *Soviet*, *anti-apartheid*) or that have increased in frequency (e.g. *clean-energy*, *anti-doping*, *low-carb*, *anti-terror*, *high def*). We can compare concordance lines for a given word or phrase, where there have been changes over time (such as *web* or *green*). As with COHA, we can compare collocates in two recent time periods. For example, collocates of *web* that were more common in the early 1990s than the late 2000s include *strand*, *spider*, *relation*, *life*, and *image*, where those from the late 2000s include *site*, *page*, *e-mail*, *browser*, and *link*. Finally, we can –

with just one quick search – see the relative frequency of related words (such as synonyms of a given word) over the past two decades.

Having now provided samples of the different types of lexical data that a large, well-designed corpus can provide, let us now turn to several examples of other corpora and lexical resources for Late Modern English, to see what types of insight they give for lexical change.

## 2. Small corpora

English historical linguistics has a strong tradition of small, well-designed corpora, in the range of one to five million words each. As López-Couso (this volume) mentions, these include the Brown family of corpora – one million words each in Brown (US 1960s), LOB (UK 1960s), Frown (US 1990s), and FLOB (UK 1990s). They also include ARCHER (1.8 million words, 1650-1999) CONCE (Corpus of Nineteenth Century Texts) (1,000,000 words, UK, 1800s), and the Helsinki Corpus (1.6 million words, Old English through the early 1700s). These are all “general” historical corpora – covering a wide range of genres and (like COHA) balanced by genre from decade to decade. There are also many small corpora of particular genres, such as letters, newspapers, or court proceedings.<sup>2</sup> Again, for a very complete and balanced overview of such corpora, please see López-Couso (this volume), especially Section 2 of that chapter.

These small corpora have certainly proven their value in research on high-frequency syntactic constructions, such as modals and other auxiliaries, pronouns, and prepositions, where even in one million words there might be hundreds or even thousands of tokens. But much less has been done – or can be done – in terms of lexical change, where there are just a handful of tokens for most words. A few studies have attempted to use these smaller corpora in looking at changes in lexis, including Leech 1992, Baker 2010, Baker 2011, Baker 2012, Baron et al 2009, Hofland and Johansson 1982, Leech and Fallon 1992, Oakes and Farrow 2007, and Sigley and Holmes 2002. But as one of the most active researchers in this field notes (Baker 2011:70):

Leech and Fallon (1992) point out that the corpora in the Brown family contain only about 50,000 word types in total, which is relatively small for lexical research, and that the majority of words will be too infrequent to give reliable guidance on British and American uses of language.

For that reason, this study focuses only on frequent words in the corpora. It was stipulated that for a word to be of interest to this study, it would need to occur at least 1,000 times when its frequencies in all four corpora were added together. Three hundred eighty words met this criteria, but a number of high frequency words (e.g., *class*, *miss*, *black*, *true*, and *English*) were excluded because they missed the cutoff.

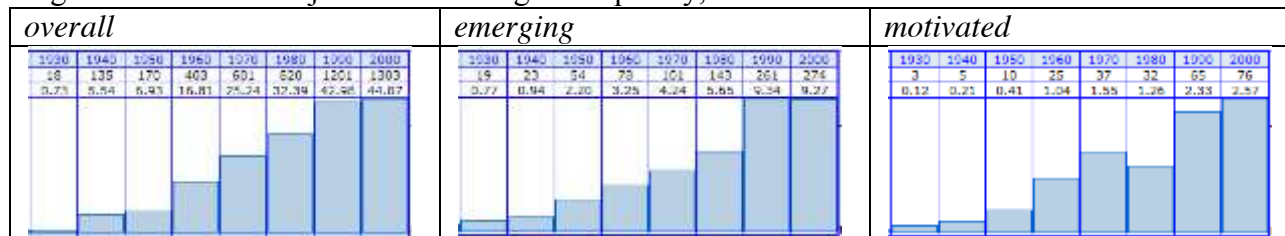
In this chapter, we will continue Baker’s line of research, and show empirically what types of lexical data we can extract from small 1-2 million word corpora, compared to a much larger corpus like COHA. As a test case, we will briefly consider adjectives that have (at least) doubled in (normalized) frequency in COHA from the 1960s to the 1990s, and then examine how well the one million word Brown and Frown corpora (US, 1960s and 1990s) provide comparable

---

<sup>2</sup> See <http://www.helsinki.fi/varieng/CoRD/index.html> for a good overview of these genre-specific corpora.

evidence for this increase in frequency. In other words, in the data below we will be considering adjectives like *overall*, *emerging*, and *motivated*, whose charts in COHA are shown below.

Figure 16. COHA: Adjectives doubling in frequency, 1960s-1990s



The following table shows that in COHA there are 15 adjectives that have a combined frequency of between 800-1600 tokens in the 1960s and 1990s (words such as *overall* (shown above), *amazing*, *long-term*, and *alternative*) and which have at least doubled in frequency during this time. There are another 127 types with a frequency of between 200-400 tokens in these two decades (e.g. *emerging* (shown above), *compelling*, *indoor*, *preferred*, and *unclear*), and 394 types with a frequency between of between 50 and 100 tokens (e.g. *motivated* (shown above), *first-time*, *blurry*, *impaired*, *viral*, *obnoxious*, and *luscious*).

Table 1. Evidence for increase in adjective frequency, COHA and Brown family

COHA: token range		800-1600	200-400	50-100
COHA: # of types		15	127	394
# Brown/Frown tokens	0	0	8	114
	1-9	1	46	264
	>= 10	Support	50	12
	>= 10	???	15	0
	>= 10	Contradict	3	4
Brown/Frown "correct"		0.40	0.39	0.03

Table 1 shows that for the 15 COHA adjectives that have at least doubled in frequency and which have a combined token frequency of 800-1600 in COHA in the 1960s and 1990s, all of these occur at least once in Brown/Frown, which is encouraging. One word occurs 1-9 times in Brown/Frown, and the other 14 occur at least 10 times (e.g. 3 tokens in Brown and 7 tokens in Frown), which is perhaps enough to show an increase from the 1960s to the 1990s. Of these 14 adjectives that occur at least 10 times, 6 do show frequency that has doubled from the 1960s-1990s (e.g. Brown 6, Frown 12, which is shown as "Support" (COHA) in Table 1 above). Another 5 adjectives show an increase, but less than the doubling in COHA (e.g. 6 Brown and 7 Frown, shown as "? ? ?" above). And in 3 cases, the Brown/Frown data actually shows a decrease from the 1960s to the 1990s (e.g. 7 Brown, 4 Frown; shown as "Contradict" above). Overall, then, 6 of the 15 types (40%) of these high-frequency adjectives in Brown / Frown show the same doubling in frequency that is shown in the robust data (800-1600 tokens) in COHA.

The situation is a bit less encouraging for the 127 medium-frequency adjectives (token count of 200-400 for the 1960s/1990s in COHA). Of these, 8 do not occur at all in Brown/Frown and 46



occur just 1-9 times, which is probably too few to see an increase. Of those occurring 10 times or more in Brown/Frown, 50 show a doubling, 15 show a smaller increase, and 8 show a decrease.

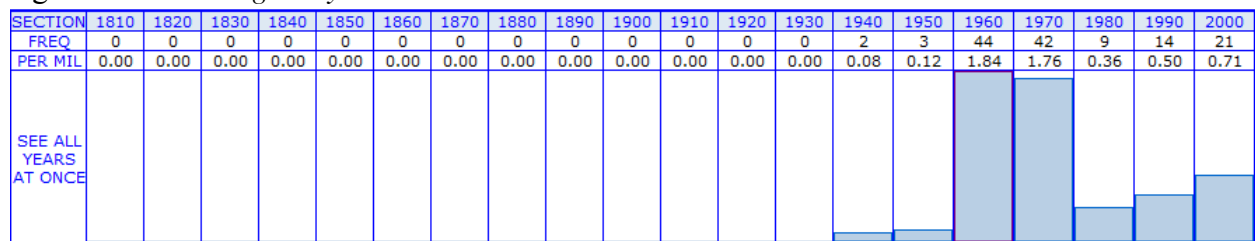
The situation with lower-frequency words is very poor. Remember, these are adjectives like *first-time*, *blurry*, *impaired*, *viral*, *obnoxious*, *luscious*, and *motivated* – less common to be sure, but certainly still the type of adjectives that most speakers of English would be familiar with. Of the 394 types in COHA with a frequency of between 50-100 and which have at least doubled in frequency, 114 of these do not occur at all in Brown/Frown, and another 264 occur less than 10 times – probably too few to be useful. As a result, Brown/Brown provides evidence for doubling in frequency for only about 3% of all of these lower-frequency adjectives from COHA.

We should also realize that for some types of searches, the situation is even much worse than the searches just described, where we are simply looking at the frequency of a given word or phrase over time. For example, Figures 11-13 above shows that collocates are extremely sensitive to corpus size. In a one million word corpus (1/400th the size of COHA), virtually none of the collocates of *gay* or *woman* would occur more than one or two times. In summary, we argue that a corpus of 1 million words – while perhaps useful for high frequency grammatical changes – is simply too small to examine lexical changes with the vast majority of the words in the language.

In addition to size, one other problem with some of these small corpora is the issue of granularity. For example, the Brown family of corpora have texts from 1961 and 1991 (and work is proceeding on a similar corpus from 1931 and then 1901). But because there are texts from only every 30 years, any changes that take place in between these years is essentially “invisible”, and in terms of lexical change, this is often too long of a gap.

Let us briefly consider two examples related to granularity, which are representative of tens of thousands of words. First, let us consider the frequency for *groovy* in COHA:

Figure 17. COHA: *groovy*



Imagine that our two corpora contained texts thirty years apart – from 1955 and 1985. In this case, it would appear (based on the COHA data from the 1950s and the 1980s) that *groovy* is on the increase. While it has increased slightly in these 30 years, we would miss entirely the steep decrease from the 1960s/1970s to the 1980s. Second, consider the case of *normalcy*:

Figure 18. COHA: *normalcy*

SECTION	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000
FREQ	0	0	0	0	0	0	0	0	0	0	0	24	27	29	25	26	27	28	34	45
PER MIL	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.94	1.10	1.19	1.02	1.08	1.13	1.11	1.22	1.52
SEE ALL YEARS AT ONCE																				

This word was famously “rescued” from obscurity by President Warren G. Harding in 1920, who (according to purists) mistakenly used it instead of the more “correct” *normality*. The word caught on with a public tired of World War I and other foreign involvements, and Harding went on to win the election. But imagine that we only had two corpora from 1901 and 1931 (as with the planned extensions in the Brown family of corpora). There would obviously be a large increase in frequency between 1901 and 1931, but there would be no way to know if that predated Harding, whether his campaign caused the increase in usage, or whether it was after his time. In summary, corpora that have texts that are spaced decades apart may be adequate for looking at much more gradual grammatical change, but they are much more problematic in looking at lexical change, which can occur quite suddenly.

### 3. Text archives

If small corpora are of limited use in looking at lexical changes, then perhaps we might turn to text archives, some of which contain hundreds of millions of words of text (see López-Couso, this volume (especially Section 3) for a very good overview of available text archives). As we will see in this section, however, while the *textual* component of these text archives is truly impressive, they are often paired with an architecture and interface that make it very difficult to extract the lexical data in any meaningful way.

In this section, we will consider how we can find lexical data using Early English Books Online (EEBO), Literature Online (LION), and the New York Times newspaper archive. Note that as we discuss these archives, the precise type of data that we can get from each text archive is not the important thing – undoubtedly this will change over the next 5-10 years, and most of this will be outdated by then. In addition, there will undoubtedly be bigger and better text archives. Rather, this discussion and these resources are meant to serve simply as *examples* of some of the challenges inherent in using text archives.

EEBO represents one end of the spectrum: text archives that – while they have excellent textual data – have an interface that makes them quite problematic for lexical research. As of late 2012, EEBO contains about 700 million words in about 28,000 texts from 1470-1699, which is quite impressive. But consider what type of information is (and is not) available from the search interface. Suppose that we search for the word *polite*:

Figure 19. EEBO entries

Results: 1715 matches in 953 records

[Refine this search](#)

1 to 25 of 953 records  
1 | 26 | 76 | 176 | 376 | 776 | 929 | [Next>>](#) Sort by: not sorted Sort Results

**Title:** The art of knowing one-self, or, An enquiry into the sources of morality written originally in French, by the Reverend Dr. Abbadie.  
**Author:** Abbadie, Jacques, 1654-1727.  
**Publication Info:** Oxford : Printed by Leonard Lichfield, for Henry Clements, and John Howell, booksellers, 1695.  
**Collection:** Early English Books Online  
— Search Results: 1 match in full text  
[Results details](#) | [Table of contents](#) | [Add to bookbag](#)

As with a typical Google web search, the results tell us how many texts the word or phrase appears in (953 texts, in this case). But there is no indication of frequency over time, as we can so easily find in COHA. Theoretically, it is possible to write a script to automatically retrieve all pages and then parse the page to extract dates and number of tokens and then import and process this in a database, but this is probably beyond the ability of most users. As far as collocates or concordances, we would need to click individually – one after another – on all 953 texts to see the word or phrase in context, as in the following.

Figure 20. EEBO concordance lines

Title: The art of knowing one-self, or, An enquiry into the sources of morality written originally in French, by the Reverend Dr. Abbadie.  
Author: Abbadie, Jacques, 1654-1727.  
Publication Info: Oxford : Printed by Leonard Lichfield, for Henry Clements, and John Howell, booksellers, 1695.  
Collection: Early English Books Online  
— Search Results: 1 match in full text  
[Table of contents](#) | [Add to bookbag](#)

The ART of KNOWING ONE-SELF: Or, An Enquiry after the Sources OF MORALITY. > The FIRST PART. Wherein we Treat of the Nature of MAN, of his End, his Perfections, his Duties, and his Strength. > CHAP. VIII. Where we continue to shew, what Effect the Sentiment of our Immortality can work upon our Heart.

• ... s too well, But that we know not our selves enough. / Such is Bashfulness, the most polite and reasonable of all the Vertues: Or rather the Artificial Disguisement of our Intemperance / ...

An example of a somewhat more useful text archive is Literature Online, which contains about 600-700 million words of text from Early Modern English and Late Modern English, and which continues to grow in size. It provides some very useful features that EEBO does not. First, it is possible to search for spelling variants and for variant forms, as well as with simple proximity and Boolean operators:

Figure 21. Literature Online (LION)

Search in:  ALL  Poetry  Drama  Prose

Search including: Variant spellings  (What is this?) (e.g. jealousy finds iealousy, jalousie etc.) Variant forms  (What is this?) (e.g. arrest finds arrested, arrests, arresting, etc.)

Keyword(s) in Work:  [select from a list >>](#)  
[check for variants >>](#)  
e.g. handful of dust; handful NEAR dust; dust\*

While it is still not possible to see the frequency of a word over time, with LION one could potentially search for a word in (for example) 1500-1509, and then 1510-1519, and so on, and after each search copy the number of hits to a program like Excel, and then generate a frequency chart from all of these searches. Finally, in Literature Online (unlike in EEBO) it is not necessary

to click on each text to see the word in context. Each page contains the word or phrase in context for up to ten texts at a time:

Figure 22. Literature Online concordance lines

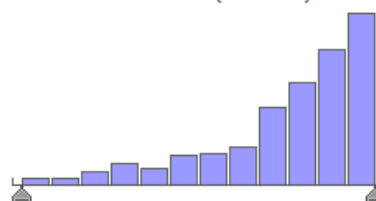
25. Anon., ca. 1603 (Philotus) [[Author Page](#)]  
[Ane verie excellent and delectabill Treatise intitult PHILOTVS. 84Kb](#) , [from Anon., ca. 1603  
 Ane verie excellent and delectabill treatise intitult Philotus (1603) ]  
[\[Durable URL for this text\]](#)  
 Found 6 hit(s).  
 ...Philotus is the man a **faith**, Ane ground-riche man and full...  
 ...zour geir I cair not: **Faith** ze zour self sall modifie,...  
 ...Quhy not? gif that with **faith** we pray For oft the...  
 ...to supplie. Pleasant. 111 Ane **faith** perfumit with fyne folie, And...  
 ...In forme of hir a **faith** I sie, Sum Deuill hes...

[View all hits in this text](#)

It would still be necessary to write a script to extract all of this information from (perhaps) hundreds of pages, and then process it in another program to create concordance lines and to extract collocates for all tokens, but it is probably still more doable than in EEBO.

As a final example of text archives, let us consider online historical newspapers. In the past 5-10 years, many historical archives for many newspapers have become available online, and many more will undoubtedly become available in the future as well. The question is whether and how this data can be used for lexically-oriented research. As an example, we will briefly consider the New York Times archive, which contains 300-400 million words from the 1850s-2000s. A search for the phrased *messed up* yields results like the following:

Figure 23. New York Times: concordance lines and frequency chart

<p>1 <a href="#">Sports of The Times: The Dodger Touch</a>          By ARTHUR DALEY. <b>New York Times (1923-Current file)</b> [New York, N.Y] 22 May 1958:          39.          ...popped into the wrong brier patch and <b>messed up</b> Herb Elliott's sub-record mile</p> <p>2 <a href="#">Sports of The Times: The Boy From Syracuse</a>          By ARTHUR DALEY. <b>New York Times (1923-Current file)</b> [New York, N.Y] 23 Jan 1958:          34.</p>	<p>3 1894 - 2009 (decades)</p>  <table border="1" style="display: none;"> <caption>Frequency of 'messed up' by decade (approximate values)</caption> <thead> <tr> <th>Decade</th> <th>Frequency</th> </tr> </thead> <tbody> <tr><td>1890s</td><td>1</td></tr> <tr><td>1900s</td><td>2</td></tr> <tr><td>1910s</td><td>3</td></tr> <tr><td>1920s</td><td>4</td></tr> <tr><td>1930s</td><td>5</td></tr> <tr><td>1940s</td><td>6</td></tr> <tr><td>1950s</td><td>8</td></tr> <tr><td>1960s</td><td>12</td></tr> <tr><td>1970s</td><td>18</td></tr> <tr><td>1980s</td><td>25</td></tr> <tr><td>1990s</td><td>35</td></tr> <tr><td>2000s</td><td>55</td></tr> </tbody> </table>	Decade	Frequency	1890s	1	1900s	2	1910s	3	1920s	4	1930s	5	1940s	6	1950s	8	1960s	12	1970s	18	1980s	25	1990s	35	2000s	55
Decade	Frequency																										
1890s	1																										
1900s	2																										
1910s	3																										
1920s	4																										
1930s	5																										
1940s	6																										
1950s	8																										
1960s	12																										
1970s	18																										
1980s	25																										
1990s	35																										
2000s	55																										

As with other text archives, some results (as in entry #1 above) yield short Keyword in Context entries, which with some effort could be aggregated to created KWIC and/or collocates for all entries. Unfortunately, many other entries (such as in #2) only provide links to a PDF file, which means that users have to scan through that entire graphics image to find the phrase *messed up* (the PDF file is not searchable), and researchers would have to do this for each of the hundreds or thousands of entries. (This type of linking to graphical/non-text PDF images is common in text archives.) In this sense, the New York Times archive is even less linguistics-friendly than Literature Online.

However, note the frequency chart (#3 above), which shows the frequency of *messed up* in each decade from the 1890s to the 2000s. Users can click on any decade to see the frequency by year, and they can also limit the KWIC entries to just that decade or year. While it is interesting and encouraging to see such linguistics-friendly features in the interface, there are some important limitations. Most importantly, the frequency data is not normalized for corpus size in each



decade. This means that researchers would have to first calculate the “base frequency” for a very common word like *the* or *and* (which presumably change very little from decade to decade), and then compare the frequency of any given word or phrase (like *messed up*) to this base frequency. Nevertheless, the inclusion of frequency information in a non-linguistic-oriented resource– in any form – is still a start.

#### 4. The Oxford English Dictionary

As anyone who has used it can testify, in many respects the OED is unparalleled in terms of the information that it can provide on lexical changes in English. As the following screen shot of the online OED interface shows, we can perform searches using an incredible range of information in the OED dictionary entries– subject, language of origin, region, usage, and part of speech. Figure 24. Oxford English Dictionary (OED) interface

The screenshot shows the Oxford English Dictionary (OED) search interface. It features a search bar at the top with a dropdown menu set to "Full Text". Below the search bar, there are two input fields for "And" and "in", both set to "Full Text", with a "Remove row" button. There are also checkboxes for "Case-sensitive" and "Exact characters". On the right side, there are "Options for NEAR/NOT NEAR" with a dropdown set to "One Word" and a checkbox for "Ordered". The main search area is divided into two columns. The left column contains filters for "Subject", "Language of Origin", "Region", and "Usage", each with a text input field and a "Browse" link. The right column contains filters for "Date of entry", "Include entries marked as", "Part of speech", and "Restrict to entry letter or range", each with a text input field and radio button options. The "Date of entry" filter includes a "Enter year or range of years" input. The "Include entries marked as" filter has radio buttons for "All", "Current", and "Obsolete". The "Part of speech" filter has a dropdown menu set to "All". The "Restrict to entry letter or range" filter has a "Enter range" input field.

For example, we could search for Subject = Arts, Language of Origin = Italian, Date of entry = 1700-1720, and Part of speech = noun, and we would retrieve the words and phrases *al fresco*, *allegretto*, *capitolo*, *claro obscuro*, *mandolin*, *mezzanine*, *pianissimo*, *ricercata*, *ridotto*, *rondo*, *smalto*, and *stuccature*.

For the purpose of this chapter, however, which is more oriented towards corpus-based approaches to lexical changes in English, we might consider how the OED can be used as a corpus – a searchable collection of texts – rather than just a collection of entries in a dictionary. Although some may not be accustomed to thinking of the OED as a corpus, we should remember

that it is based on approximately 38 million words of text from about 2.8 million entries, which means that we could – at least in principle – search it as a “corpus”.

At the most basic level, we can search for a word or phrase and see that word or phrase in context. Consider first a few entries for the word *polite*:

Figure 25. OED: quotations by headword

<p><b>a. Smoothed, polished, burnished. <i>Obs.</i></b></p> <p>1398 J. TREVISA tr. Bartholomaeus Anglicus <i>De Proprietatibus Rerum</i> (BL Add.) f. 197<sup>r</sup>, Berill is..yliche in grene colour to Smaragde but it is wiþ palenesse and polit.</p> <p>1429 <i>Mirour Mans Saluacioune</i> (1986) l. 1499 The Arche withinne &amp; without was hiled with golde polyt.</p> <p>1488 (1478) HARY <i>Actis &amp; Deidis Schir William Wallace</i> (Adv.) x. l. 388 Throu polyt platis with poyntis persyt thair.</p>	<p>Thesaurus »</p>
<p><b>c. Courteous, behaving in a manner that is respectful or considerate of others; well-mannered.</b></p> <p>1751 E. JUSTICE <i>Amelia</i> 193, I am much concern'd that the Gentleman should make you pay Ten Shillings for that trifling Present I sent you, but that Family is not capable of a polite Action.</p> <p>1772 H. MACKENZIE <i>Man of World</i> (1823) II. xx. 492 The French are the politest enemies in the world.</p> <p>1781 GIBBON <i>Decline &amp; Fall</i> II. xix. 151 Narses..was endowed with the most polite and amiable manners.</p>	<p>Thesaurus »</p>

In the case of the COHA data shown above (Figures 2-3, 11), we used re-sortable concordance lines and collocates to separate entries with different meanings. In the case of the OED, this has already been done for us, as seen above for the two meanings of *polite*.

But the OED interface also allows us to search for *any* words in the 2.8 million quotations, whether or not they are the headword (e.g. *polite*, above). For example, we could search for *trifling* and see the quotation for the headword *polite* from 1751 above or search for *endowed* and see the quotation for the headword *polite* from 1781. As a concrete example, the following are a handful of entries for *trifling*:

Figure 26. OED: searching for quotations by keyword

<p><b>11. <i>polite, adj. and n.</i></b></p> <p>...you pay Ten Shillings for that trifling Present I sent you, but that...</p>	<p>1751</p>
<p><b>12. <i>practiser   practicer, n.</i></b></p> <p>...kes of Practisers, and Want of trifling Forms, may nonsuit you?...</p>	<p>1741</p>
<p><b>13. <i>puerile, adj. and n.</i></b></p> <p>...ery Remarks Swift 1752 78 They are trifling and I had almost said pueril...</p>	<p>1751</p>

Even with such entries, however, we would still need to write a script to parse the web pages and extract the text and the dates from all of the quotation. In addition, as displayed in this view

(Advanced / Quotations search), the quotation only show 5-6 words to the left and to the right of the search term (rather than the entire quotation), which is a bit limiting. Finally, from a corpus linguistics perspective (which is the focus of this chapter), we might also wish that we had somewhat more control over the quotation searches in the OED, by being able to search in more advanced ways for strings of words (including variant forms and substrings) and then being able to quickly and easily display the results in frequency charts, concordance lines, or with collocates, none of which is possible.

The fact that the OED is still limited as a linguistic *corpus* should not, however, take away from its unique role as a resource for lexical changes in English, through the massive amount of data in its dictionary entries. And hopefully, at some point, the interface for the OED will be improved (as it has been in many other respects over the past few years), to allow an even greater range of lexically-oriented searches of the 38 million words of data in the 2.8 million quotations.

## 5. Google Books

The Google Books n-grams are based on hundreds of *billions* of words in more than a million books. It is composed of a number of different “datasets”, such as American English (155 billion words), British English (34 billion words), and the “One Million Books” dataset (89 billion words). When the Google Books n-grams were released in late 2010, it was heralded as a resource that would revolutionize not only historical linguistics, but also a wide range of other fields within the scope of what was called “Culturomics”, such as history, cultural studies, and sociology, by showing trends in word and phrase usage over the past few centuries (see Michel, Lieberman, et al, 2010).

The n-grams – single words and phrases up to five words in length – are searchable via the web interface, and they show the frequency of the words and phrases in each decade and year. For example, the following three charts show the frequency from the 1810s-2000s for the words *teenager*, *steamship*, and *many a time*. Compare these to the equivalent searches in COHA (Figures 1-3 above), and note how similar the frequency data is in the two resources.

Figure 27. Google Books: *teenager*

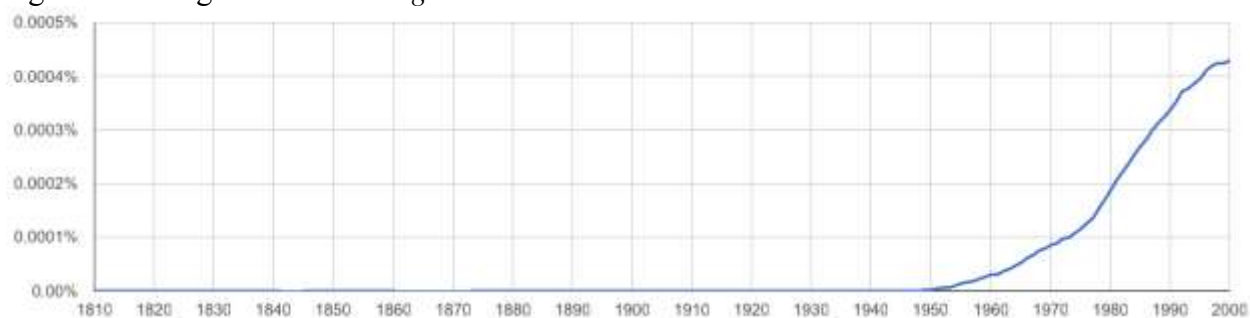


Figure 28. Google Books: *steamship*

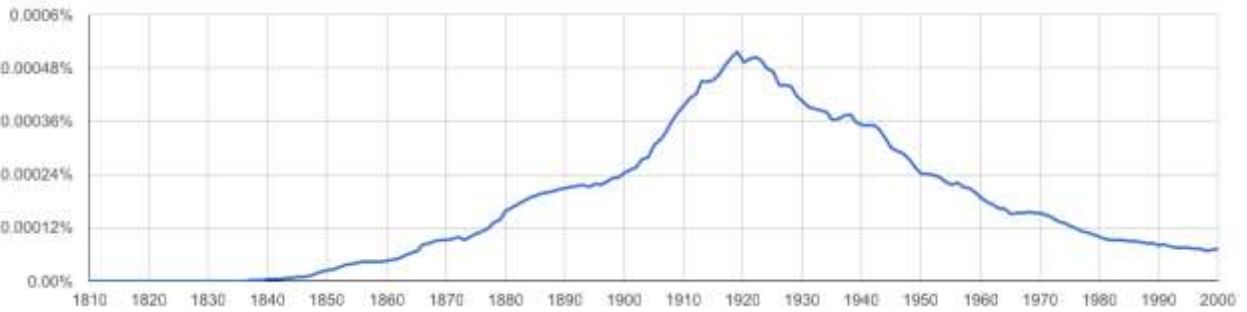


Figure 29. Google Books: *many a time*



While the total number of texts and words for the Google Books n-grams is huge, there are some critical weaknesses in terms of the interface to this data. First, the only frequency data are the pictures themselves – there is no actual frequency data that can be copied and pasted from the Google Books website. For example, [ 0.0002 ] in the chart for *many a time* is not an actual number – it is just a picture of that number. Second, the frequency data does not show the number of tokens – just these percentages. There is no way of knowing how many times a word or phrase actually occurs in a given decade or year. Third, because there are no actual token counts, it is impossible to calculate the total number of tokens for a given construction or even a single lemma. For example, one can search for *steamship* and then for *steamships* (cf. Figure 28), but these two “pictures” cannot be “added” together, as if there were actual frequency data. Fourth, it is not possible to search by part of speech. For example, one can search for *many a time* (Figure 29), *many a man*, or *many a mile*, but there is no way to search for *many a* [noun].

Finally, the only information available is the frequency of a given word or phrase over time. It would be difficult or impossible to create concordance lines or collocates for a given word or phrase. Researchers can click on a range of years in the chart to see Google-like “snippet” views with the word in context, as in Figure 30. But depending on the copyright status of the book, the context view may (#1 below) or may not (#2) be available. Even if it is, one would have to write a script to process the snippet entries and extract the relevant text (and years), and this is made difficult by the fact that Google limits snippet views to the first 1000 entries.

Figure 30. Google Books: “snippet” view

1	<p><a href="#">Improving teenage nutrition - Page 8</a>  <a href="https://books.google.com/books?id=p2gvAAAAAYAAJ">books.google.com/books?id=p2gvAAAAAYAAJ</a>            United States. Dept. of Agriculture. United States. Federal Extension Service - 1963 - Read            Friends. The <b>teenager</b> wants his peer group to like him. ... The <b>teenager</b> wants to be with others of his age group. ... The <b>teenager</b> needs to know how nutrition contributes to shiny hair, clear skin, good posture, and, of course, correct weight.</p>
---	--



2	<a href="https://books.google.com/books?id=5flC AAAAIAAJ">Adolescence and religion: the Jewish teenager in American society</a> <a href="https://books.google.com/books?id=5flC AAAAIAAJ">books.google.com/books?id=5flC AAAAIAAJ</a> Bernard Carl Rosen - 1965 - Snippet view
---	--

## 6. Summary of resources

At this point, it might be useful to review what types of lexical data can be obtained from different types of resources<sup>3</sup>. The following table might serve as a useful summary. Again, we recognize that the features of a given resource will change over time, and that new resources will become available. The point here is to see how well even the best resources are at providing information on lexical change.

Table 2. Lexically-oriented queries with different resources

	COHA	Small	EEBO	LION	NYT	OED	GB
1. Concordances for a particular word or phrase, to see the patterns in which it occurs	+	((~))	((+))	((+))	((+))	((+))	((+))
2. The frequency over time of specific words and phrases	+	((~))			(+)		+
3. Lists showing all words that are more frequent in one period than in another	+					(+)	
4. Collocates of a particular word or phrase (over time) in order to see semantic change	+						
5. Easy access to the relative frequency of related words, such as synonyms or words in a user-defined list	+	((~))				((+))	

[1] Other than COHA, none of the resources does a particularly good of providing concordance lines for all tokens of a word or phrase, although with some scripting and programming (or third party software, like WordSmith or AntConc) this can be extracted. [2] COHA and Google Books provide frequency data over time, but Google Books is limited by just showing pictures of frequency, with no real data, and the New York times frequency data is not normalized. [3] The OED comparisons of all words are based on first occurrences of words, but (unlike COHA) not on actual frequency data. [4] Other than COHA, none of the resources provide information on collocates. [5] The OED shows related words, via the Historical Thesaurus of English, but unlike COHA, it does not provide data to see the relative frequency of these words over time. Finally, a few of the features are possible with small corpora (Section 3), but only for high frequency words – and only if they are somehow accessed with the right corpus interface.

<sup>3</sup> COHA = Corpus of Historical American English, EEBO = Early English Books Online, LION = Literature Online, NYT = New York Times, OED = Oxford English Dictionary, GB = Google Books

## 7. Extending useful lexical searches to other corpora

As we have seen, COHA (and COCA) are able to carry out a wide range of searches on lexical changes in Late Modern English, which in many cases would be difficult or impossible with other corpora. However, it is important to understand that once this architecture and interface are available for corpora like COHA and COCA, they can in principle be applied to any other set of texts, including those from text archives. For example, the TIME Magazine corpus (Davies 2012c) contains 100 million words in more than 275,000 texts from the Time Magazine text archive from the 1920s-2000s (<http://www.time.com/time/archive>), and all of the lexically-oriented searches that are possible with COHA and COCA can be carried out in this “corpus” as well.

In addition, a modified version of the COHA/COCA interface has been applied to the Google Books n-grams data (see <http://googlebooks.byu.edu>), to allow a wide range of searches on this 200+ billion words of data, which are not possible with the standard Google Books interface. In the interest of space, we will provide only four examples of the types of searches that can be done with the BYU/Advanced Google Books interface (GB-BYU), which are not possible with the standard Google Books interface (GB-Standard), but many more sample searches are available at the corpus website.

First, the GB-BYU interface shows the raw and normalized frequency for words and phrases – they are not just pictures of frequency as in GB-Standard:

Figure 31. Frequency of *steamship* in BYU/Advanced Google Books (GB-BYU)



Because there are actual numbers, the charts in GB-BYU can represent the totals for lemmas or entire constructions as well. It is also possible to see the frequency of each matching string, as with the construction *many a* [noun] (cf. COHA in Figure 8 above). Note that GB-BYU allows searching by part of speech, unlike GB-Standard.

Figure 32. GB-BYU: Frequency of *many a* [noun]

WORD(S)	CHARTS	TOTAL	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000
1 many a time	G B	100293	238	534	1535	2318	4440	3915	4832	7827	11123	11791	11680	7343	1057	3721	4081	5483	4141	2997	3341	2881
2 many a day	G B	64963	131	423	1084	1922	3555	2598	3461	5173	8067	8012	7975	4449	2874	1088	2265	3198	2358	1452	1704	1861
3 many a man	G B	60128	117	369	790	1896	2012	1877	2381	3678	5204	6773	8161	8224	3360	3605	1004	4427	2748	1813	2134	2201
4 many a year	G B	53485	182	514	924	1532	2789	2178	2726	4182	6197	5768	6448	3789	2006	2083	2148	2839	2010	1385	1587	1880
5 many a night	G B	15933	38	86	208	395	686	576	623	932	1371	1488	1746	1117	687	660	734	841	813	735	924	1338
6 many a mile	G B	12920	40	76	291	473	809	586	724	1246	1392	1441	1817	898	399	430	444	855	437	289	303	269
7 many a heart	G B	12174	49	146	462	965	1577	1007	895	1174	1587	1193	1116	479	333	219	338	367	381	194	218	235
8 many a woman	G B	8821	20	26	66	123	248	229	265	488	732	944	1130	702	499	374	474	516	609	387	497	535
9 many a battle	G B	7074	13	27	110	240	381	499	380	395	895	746	752	494	320	323	320	450	271	188	223	278
10 many a tale	G B	6597	21	61	239	328	529	396	372	480	613	642	771	378	311	210	197	303	234	181	165	273

As with COHA, but unlike GB-Standard, the GB-BYU interface and architecture allows users to use wildcards to see the frequency of all matching forms in each decade, for example the

frequency of *\*ism* words (note the increase in *criticism*, *organism*, and *capitalism*, and the decrease in *baptism* and *patriotism*):

Figure 33. GB-BYU: *\*ism* words by decade

WORD/ID	CHARTS	TOTAL	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960
1 criticism	G B	5426220	1120	7951	17180	26706	42620	79076	12788	100956	169241	247306	346591	200083	245163	240566	333238	618744
2 mechanism	G B	4885909	1327	2517	7553	9308	15781	11640	19417	34598	53826	84166	177370	145663	190132	188934	271249	501658
3 organism	G B	2756106	25	70	395	4646	14857	16213	31312	56853	97622	159625	273762	177218	136206	135460	205369	313376
4 metabolism	G B	2071179	5		5			2	89	2054	7616	40725	95471	67019	53484	61663	105854	199683
5 Judaism	G B	1460561	910	1163	3461	7539	4545	9986	14113	21564	46153	36752	30449	30322	36387	46835	81562	131255
6 capitalism	G B	1427387	1		2	18	5	3	9	522	1628	7820	22454	29741	74921	75324	72496	155756
7 baptism	G B	1369446	19435	16361	49640	85027	91912	51310	84473	70753	73922	85679	80490	38122	26353	34339	51327	87642
8 socialism	G B	1181815	5	3	3	258	917	666	1262	7143	21972	33349	55803	44502	50594	61740	82035	187268
9 patriotism	G B	1175726	5406	10802	21882	31023	48202	42091	34389	54712	74083	93479	137519	83007	56389	48637	49261	84120

Like COHA, but unlike GB-Standard, in GB-BYU one can find all words that are more common in one period than another, for example *\*ism* words more common in the 1860s-1910s or the 1970s-2000s (compare Figure 10 from COHA, above):

Figure 34. GB-BYU: Comparison of *\*ism* words, 1860s-1910s vs 1970s-2000s

SEC 1: 32.8 BILLION WORDS (1860-1919)						SEC 2: 76.2 BILLION WORDS (1970-2009)					
WORD/PHRASE	1: 1860-1919	2: 1970-2009	#/BL 1	#/BL 2	RATIO	WORD/PHRASE	1: 1860-1919	2: 1970-2009	#/BL 1	#/BL 2	RATIO
1 metamorphism	1,834	22	55.8	0.3	193.46	1 consumerism	66,941	1	1,140.7	0.0	37,484.05
2 apogeeotropism	1,464	25	44.6	0.3	135.90	2 existentialism	66,111	1	867.4	0.0	28,488.12
3 hemihedism	1,274	46	38.8	0.6	64.27	3 environmentalism	47,385	1	621.7	0.0	20,416.84
4 catheterism	3,700	151	112.7	2.0	56.86	4 Surrealism	46,800	1	614.0	0.0	20,166.75
5 katabolism	6,366	265	193.8	3.3	55.75	5 isolationism	42,459	1	557.1	0.0	18,296.16
6 diorotism	1,649	69	50.2	0.9	55.46	6 Racism	161,705	5	2,121.6	0.2	13,936.17
7 Anarism	11,099	498	337.9	6.5	51.72	7 racism	818,513	27	10,738.9	0.8	13,063.27
8 aneurism	75,991	4,072	2,313.7	53.4	43.31	8 Sexism	46,226	2	606.5	0.1	9,959.70
9 Traumatism	4,001	238	121.8	3.1	39.01	9 McCarthyism	35,259	2	462.6	0.1	7,596.79
10 dilettanteism	2,569	162	78.2	2.1	36.80	10 minimalism	17,005	1	223.1	0.0	7,327.68
11 sinapism	1,423	95	43.3	1.2	34.76	11 sexism	193,193	12	2,534.7	0.4	6,937.46
12 Hindooism	1,941	130	59.1	1.7	34.85	12 Pentecostalism	24,987	2	327.8	0.1	5,383.62

GB-BYU allows users to use semantically-oriented queries with synonyms, such as [=beautiful] *woman* (compare the COHA data in Figure 15 above):

Figure 35. GB-BYU: Synonyms: “beautiful” + *woman*

WORD/ID	CHARTS	TOTAL	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000
1 beautiful woman	G B	180051	218	485	1206	1670	2889	2368	3432	5535	8643	11134	11736	8330	6027	5458	7429	11403	12181	15608	25118	37801
2 attractive woman	G B	34629	13	15	18	20	22	26	32	39	47	57	68	81	93	106	146	247	395	595	858	1057
3 lovely woman	G B	33754	120	139	576	780	1198	1333	1195	2146	2532	2733	2416	1801	1481	1155	1299	1920	1814	2097	2970	4259
4 handsome woman	G B	29458	76	116	321	433	750	617	798	1394	2254	2608	2566	1791	1400	1304	1431	2206	2004	2087	2478	2821
5 charming woman	G B	32237	43	123	399	272	451	333	420	967	3042	2981	2502	1900	1434	1011	1101	1479	1246	1000	1194	1250
6 wonderful woman	G B	18959	9	19	52	85	85	178	214	269	347	1042	1615	1403	855	843	956	1188	1211	1665	2570	3658
7 striking woman	G B	2375	1		1		1	1	1	18	14	27	23	25	31	46	67	107	188	326	681	1017
8 delightful woman	G B	2465		5	19	33	14	32	32	65	137	220	285	245	143	93	133	177	144	170	233	321
9 gorgeous woman	G B	2451						1	5	4	5	13	3	34	31	22	44	92	347	496	1432	
10 magnificent woman	G B	2064		1	8	25	44	48	78	91	111	161	136	114	81	73	79	146	144	182	244	361
11 stunning woman	G B	1337						1	4	1	7	27	47	12	20	20	21	38	70	138	295	655

Like COHA, but unlike GB-Standard, in GB-BYU one can compare collocates in different periods, for example adjectives with *woman* in the 1850s-1910s (left) and the 1970s-2000s (right) (compare the COHA data in Figure 13 above):

Figure 36. GB-BYU: Comparison of collocates with *women*, 1850s-1910s vs 1970s-2000s

SEC 1: 35.8 BILLION WORDS (1850-1919)						SEC 2: 76.2 BILLION WORDS (1970-2009)					
WORD/PHRASE	1: 1850-1919	2: 1970-2009	P/BOL 1	P/BOL 2	RATIO	WORD/PHRASE	1: 1850-1919	P/BOL 1	P/BOL 2	RATIO	
1 feeble women	816	149	22.8	2.0	11.66	1 bisexual women	10,784	1	141.5	0.0	5,064.71
2 fair women	8,212	1,529	229.4	20.1	11.44	2 battered women	83,346	10	1,093.5	0.3	3,914.35
3 Fair women	524	104	14.6	1.4	10.73	3 heterosexual women	21,388	4	280.6	0.1	2,511.22
4 chief women	699	139	19.5	1.8	10.71	4 academic women	4,253	1	55.8	0.0	1,997.42
5 delicate women	2,651	605	74.1	7.9	9.33	5 negative women	3,696	2	48.5	0.1	867.91
6 defenceless women	1,474	341	41.2	4.5	9.20	6 urban women	10,521	7	138.0	0.2	705.88
7 tender women	816	190	22.8	2.5	9.14	7 Black women	102,960	89	1,250.8	1.9	709.80
8 noblest women	975	229	27.2	3.0	9.07	8 overweight women	4,287	3	56.2	0.1	671.13
9 handloomst women	1,037	249	29.0	3.3	8.87	9 Inuit women	1,163	1	15.3	0.0	546.20
10 nervous women	1,895	476	52.9	6.2	8.48	10 Jamaican women	1,133	1	15.1	0.0	541.51

As we have seen, GB-Standard and GB-BYU use exactly the same n-grams datasets, which Google has graciously made available to others. But due to its unique architecture and interface, GB-BYU allows a wide range of lexically-oriented searches that are quite impossible with the standard, simplistic Google Books interface.

To conclude this section, one might ask why one would want to use a “small” corpus like COHA (at a “mere” 400 million words) when they could use the Google Books n-grams, which is based on a dataset that is nearly 400 times as large, at 155 billion words.

First, Google Books are just n-grams – nothing longer than a five word string can be searched (where in COHA, the search string can be up to 21 words in length). Even more importantly, these Google Book n-grams are limited to just those that occur 40 times, which dramatically reduces the number of types available (where in COHA, *all* strings – even those occurring just once or twice – are included). In addition, there is no contextual part of speech disambiguation in Google Books (as there is in COHA), making it impossible to search for “nouns” or “verbs” (or any other part of speech) in isolation. And finally, Google Books (whether GB-Standard or GB-BYU) only displays the actual text – the original sentences and perhaps paragraphs – in individual “snippets” (see Figure 30 above), whereas they are displayed in COHA in a much more readable and usable format (see Figure 1-3, 9).

So even with the important improvements that GB-BYU makes on the basic GB-Standard interface, there are still many types of lexically-oriented searches that are only available with COHA, since it alone of the two is a true “linguistic corpus”.

## 8. Early Modern English

Nearly all of the discussion to this point has focused on Late Modern English, via COHA, Google Books, text archives, and smaller corpora like the Brown family of corpora. But there is no reason that the same interface that is used for COHA, COCA, TIME, and GB-BYU could not be applied to texts from Early Modern English as well, to provide for a wide range of searches on lexical change and variation. In fact, we have applied this same interface to 400 million words of data in nearly 30,000 texts from Early English Books Online. Due to space limitations here, we will provide only a handful of examples of lexically-oriented searches from EEBO-BYU.

As with the other corpora, we can see concordance lines for a particular word, phrase, or construction. For example, the following are a handful of lines from EEBO-BYU for the word *wise* from the 1470s-1550s:

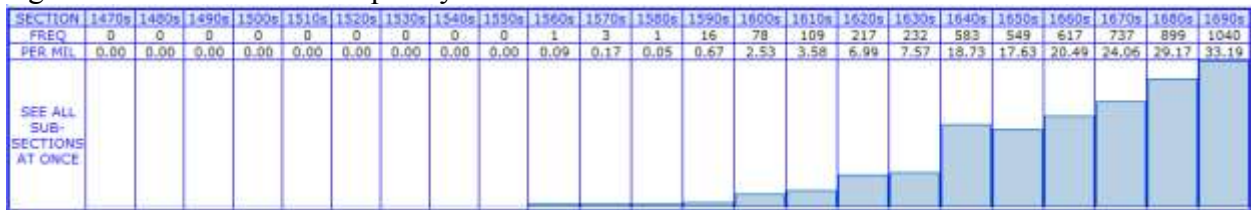


Figure 37. EEBO-BYU: Concordance lines for *wise*, 1470s-1550s

36	1506	The mirroure of gold	A	B	C	of Iugemente shall see the blessed creatures / not in that <b>maner</b> <b>wyse</b> to knowe their love what it is : but only they shall
37	1507	Scala perfectionis.	A	B	C	and purified twenty wynter / And <b>therefore</b> <b>upon</b> <b>this</b> <b>maner</b> <b>wyse</b> take in my <b>avenge</b> as I have sayd .. and na 3-ly.
38	1534	A mustre of ecismety	A	B	C	sende esury good men in a iust <b>cause</b> / <b>rather</b> <b>more</b> wyser counsaile than he <b>were</b> lyke to haue of suche men . whiche
39	1549	A contrarye (to a ce	A	B	C	; preacher of Christ . Wherefore Christe <b>willets</b> <b>not</b> <b>in</b> <b>my</b> <b>wyse</b> these kodes or <b>darnel</b> of hereti-ques , and not of Adulterers
40	1533	The debellacyon of S	A	B	C	Page xviiiAnd therefore me thynketh that this <b>deuice</b> <b>is</b> <b>not</b> <b>mych</b> wyser , than the <b>de-vice</b> that a good felow deuised ones for his
41	1553	The vocacyon of Ioha	A	B	C	my tryall / I desired to go a litle / <b>but</b> <b>in</b> <b>no</b> <b>wyse</b> <b>wolde</b> <b>it</b> <b>be</b> <b>graunted</b> . * After that we passed more than the
42	1537	The boke named the G	A	B	C	sustained at the batayle of Cannas , they <b>could</b> <b>in</b> <b>no</b> <b>wyse</b> <b>deleyer</b> them . <b>Wherefore</b> they dyscharged them of their prymyse ,

Second, we can search for frequency of words and phrases in Early Modern English. In the following, we search for the verb *secure* as part of the string *to secure*, to distinguish it from *secure* as an adjective, since the corpus is not contextually disambiguated.

Figure 38. EEBO-BYU: Frequency of *to secure*



As with the other corpora, we can see the frequency of all matching strings by decade, as with words ending in *\*ism*. (Note that earlier spellings with *\*isme* or *\*ysm(e)* are not included here.)

Figure 39. EEBO-BYU: *\*ism* words by decade

	CONTEXT	ALL	1470s	1480s	1490s	1500s	1510s	1520s	1530s	1540s	1550s	1560s	1570s	1580s	1590s	1600s	1610s	1620s	1630s	1640s	1650s	1660s	1670s	1680s	1690s
1	BAPTISM	14051								8	6	12	9	75	28	74	251	185	146	578	2885	1757	2309	2785	2943
2	SCHISM	4862												6	1	10	16	23	3	150	367	667	900	1487	1132
3	ATHEISM	1833												1	2	6	6	11	5	23	119	392	371	399	497
4	CATECHISM	1396											1	6	3		5	3	20	34	90	189	279	421	346
5	SYLLOGISM	634												3	2	1	1	8		17	45	63	123	198	257
6	PAGANISM	352													1	1	16	3		21	40	129	105	111	123
7	INFANT-BAPTISM	524																		9	173	38	62	30	154
8	ANTI-CHRISTIANISM	332														2	18			11	7	218	43	21	12
9	JUDAISM	314																		9	49	43	65	87	81
10	QUAKERISM	284																			2	12	96	33	121

In a text of Early Modern English, there are of course several challenges. The first, of course, is spelling variation (see Smitterberg, this volume). However, we have received from Martin Mueller at Northwestern University a database of more than 1,350,000 variant spellings for more than 1,020,000 word forms (from approximately 73,000 lemmas). Because these forms are in a database that is linked to the standard search interface and architecture, we can include variant forms as part of queries. For example, Figure 40 shows a few of the 25 different forms for *music* and their frequency in each decade, and Figure 41 shows the frequency for a few of the variants of *loud music*.

Figure 40. EEBO-BYU: Variant spellings of *music*, by decade

CONTEXT		ALL	1470s	1480s	1490s	1500s	1510s	1520s	1530s	1540s	1550s	1560s	1570s	1580s	1590s	1600s	1610s	1620s	1630s	1640s	1650s	1660s	1670s	1680s	1690s
1	MUSICK	7310		2					1	4	31	92	162	254	387	473	250	407	513	841	1170	934	1001	1107	
2	MUSICKE	6212						3	15	21	84	304	484	1038	1129	939	981	892	286	41	13				2
3	MUSIKE	600		6	2			1	33	25	47	64	100	79	68	108	87	23	27	2	4			4	
4	MUSIC	232		1				1	1	1	3			1		2	9		1	3	3	32	30	63	79
5	MUSICKS	219												12	5	20	53	27	16	13	17	14	6	10	8
6	MUSYKE	199		2	8	7	4	16	3	46	24	31	15	2	1										
7	MUSICKES	135												5	11	17	26	26	19	13	7			1	
8	MUSIK	70		1	9				1	4	9			14	21	2	5	1	2					1	1

Figure 41. EEBO-BYU: Variant spellings of *loud music*, by decade

CONTEXT		ALL	1470s	1480s	1490s	1500s	1510s	1520s	1530s	1540s	1550s	1560s	1570s	1580s	1590s	1600s	1610s	1620s	1630s	1640s	1650s	1660s	1670s	1680s	1690s	
1	LOUD MUSICK	87																	1		2	2	9	18	5	2
2	LOUD MUSICKE	35																	4	3	5	17	3			
3	LOWD MUSICKE	11																	4		1	6				
4	LOWDE MUSICKE	5																	4	1						
5	LOUDE MUSICKE	4																		1	2		1			
6	LOWD MUSICK	4																			2	1		1		

The second challenge is part of speech (see Smitterberg, this volume). We have extracted part of speech data from the Penn Helsinki Parsed Corpus of Early Modern English (PCEME). Words are assigned a given part of speech in the EEBO-BYU lexicon if they occur three times or more with that part of speech in the PCEME. For example, *poore*, *certaine*, *generall*, *faire*, *diuers*, *naturall*, *divers*, *lawfull*, *lyke*, *greate*, *neere*, *equall*, *litle* are all (non-PDE spelling) adjectives in our lexicon, because each of these is tagged as an adjective at least three times in the PCEME.

Although this solution is not perfect, it does allow it to use provisional part of speech tags in searches, as in the following collocates of *woman* in EEBO-BYU.<sup>4 5</sup>

Figure 42. EEBO-BYU: Part of speech: adjective collocates of *woman*

#	CONTEXT	ALL	1470s	1480s	1490s	1500s	1510s	1520s	1530s	1540s	1550s	1560s	1570s	1580s	1590s	1600s	1610s	1620s	1630s	1640s	1650s	1660s	1670s	1680s	1690s
64	WANTON	170						2	6	5	8	9	9	16	18	19	20	17	10	7	8	9	2	8	
65	CHAST	169	1	1	1		1	2	5	7	12	5	6	15	22	10	19	16	8	6	5	4	8	13	
66	FIT	189							1		1	3		11	19	25	15	20	13	12	10	10	10	16	
67	LEWD	188			1								2	2	7	3	10	21	12	18	17	24	11	14	27
68	BEAUTIFUL	188											8	2	1	7	3	1	3	2	2	14	18	37	45
69	DIUERS	168						1	3	6	12	14	21	30	27	23	24	6							
70	LEARNED	168							2	3	6	9	11	11	14	7	10	8	4		29	18	19	9	11
71	LAST	166		1				1	2	7	4	1	5	9	11	15	13	13	4	9	15	14	11	30	
72	WRETCHED	165				1				6	9	9	6	10	23	8	8	20	22	8	8	7	8	7	10

The third challenge are the synonyms. Obviously, a modern thesaurus (as we use with COHA or COHA or BYU Google Books) will be inadequate for Early Modern English. To compensate for this, we use an electronic version of the Historical Thesaurus of English (HTE), and allow users to create “customized word lists” based on the HTE. For example, the following table shows words from the HTE relating to “sweet” (note that variant spelling forms are included as well):

<sup>4</sup> Because entries are taken from further down in the list, the decade headers are not shown for the columns in this figure or in the following two figures.

<sup>5</sup> *fit* (#66) is probably not an adjective in these cases (since there is no contextual disambiguation), the other words probably are adjectives

Figure 43. EEBO-BYU: Word lists based on the *Historical Thesaurus of English*

42	EXQUISIT	325										9	9	17	35	30	32	25	52	18	26	5	15	13	17		
43	SAHORE	295					2					4	13	21	33	33	34	50	43	48	1						
44	SOFTELY	272	8	17	3	7	4	30	39	38	17	43	32	17	11	2				4							
45	UNCTUOUS	272																	11	11	52	64	95	79	39		
46	SAINTEST	269										1	5	16	38	23	42	54	25	15	18	9	8	9	4		
47	LUSTYE	252						3	15	41	33	18	18	24	6			1									
48	DELICATE	237	1	14	3	5	9	16	80	38	21	15	18	24	3			1	1								
49	PALATES	229								1			3					3	14	15	28	25	24	24	24	25	31
50	SAUDRY	210						1	4	3	3	10	13	17	24	18	38	38	29	13							
51	LUSTELY	208	2	4	1	1	9	6	15	25	9	27	34	33	22	18	3	1									

These HTE-based customized word lists can then be used as part of queries, such as the following, which finds words in the HTE-based list for “pleasant” followed by words in the HTE-based list for “music” (note the alternate spellings as well). While certainly not perfect, this approach does allow users to perform semantically-oriented searches, rather than just simply searching for exact strings.

Figure 44. EEBO-BYU: Synonym-based searches: “pleasant” + “music”

19	PLEASANT SONGES	5							1			1	2	1													
20	COMELY COMPOSITION	3												1	1				1		2						
21	DULCET TUNES	4														3					1						
22	PLEASANT MELODY	4									1	1	2														
23	PLEASANT MELODIE	4										1	1	1	1												
24	PRETTY SONGS	4																		1		1				2	
25	PLEASANT SONG	3										1	1	1													
26	DOUCE HARMONIE	3														3											
27	AGREEABLE MUSICK	3																						1		1	1
28	AMICABLE COMPOSITION	3																						1		1	

These are just a handful of examples of how we can use a robust textual corpus (the 400 million word EEBO) and a robust search interface to carry out lexically-oriented research on Early Modern English. Several problems remain, and the corpus is not yet publicly-available (due to licensing issues with the Text Creation Partnership, the creators of EEBO). Nevertheless, it will hopefully give an idea of what can be done even with Early Modern English texts, with the right tools.

## 9. Discussion and conclusions

As we have discussed in this paper, there are two crucial features of corpora that need to be present, in order for us to be able to look at this wide range of lexical (and semantic) shifts. The first is size. As we have discussed at some length in Section 2, small 1-2 million word corpora are perfectly fine for looking at high frequency syntactic constructions, but they are usually far too small to adequately investigate medium and low-frequency words and phrases. But size is not everything. As we have also discussed in Section 3-5, the data in extremely large text archives is essentially “trapped” inside, unless we have the right corpus architecture and interface to efficiently access this data.

As we have seen in Section 1 (as well as Section 7), there are a handful of corpora (such as COHA, TIME, and Google Books – Advanced) that have recently become available, which are both very large and which also have a robust architecture and interface. With these corpora, we

can (at the most basic level) use re-sortable concordances to search for individual tokens of a word or phrase, and to see the patterns in which they occur. We can also find the frequency of a word or phrase over time, and (as a more advanced search) we can produce lists of all words that are more frequent in one period than in another. Moving more towards semantic change, with these corpora we can also use collocates to look at changes in meaning over time. And since words do not occur in isolation paradigmatically, we can search for the frequency and use of related words over time.

In summary, due to recent advances in English historical corpora, we can now examine lexical change in English in a number of ways that would have been quite impossible even 4-5 years ago. And because these new architectures and interfaces can be applied to new datasets (such as Early English Books Online, as discussed in Section 8), the next few years will likely see even more advances in terms of examining historical English lexis.

## References

- Baker, Paul. 2009. The BE06 Corpus of British English and recent language change. *International Journal of Corpus Linguistics* 14(3). 312-337.
- . 2010. 'Will Ms ever be as frequent as Mr? A corpus-based comparison of gendered terms across four diachronic corpora of British English.' *Gender and Language* 4.1: 125-129.
- . 2011. 'Times may change but we'll always have money: a corpus driven examination of vocabulary change in four diachronic corpora.' *Journal of English Linguistics* 39: 65-88.
- Baron, Alistair, Paul Rayson, and Dawn Archer. 2009. Word frequency and key word statistics in corpus linguistics. *Anglistik* 20: 41-67
- Davies, Mark. 2011. "The Corpus of Contemporary American English as the First Reliable Monitor Corpus of English". *Literary and Linguistic Computing* 25: 447-65.
- . 2012a. "Expanding Horizons in Historical Linguistics with the 400 million word Corpus of Historical American English". *Corpora* 7: 121-57.
- . 2012b. "Examining Recent Changes in English: Some Methodological Issues". In *Handbook on the History of English: Rethinking Approaches to the History of English*, eds. Terttu Nevalainen and Elizabeth Closs Traugott. Oxford: Oxford Univ. Press. 263-87.
- . 2012c "Recent shifts with three nonfinite verbal complements in English: Data from the 100 million word TIME Corpus (1920s-2000s)". In *Current Change in the English Verb Phrase*, ed. Bas Aarts, et al. Cambridge: Cambridge Univ. Press. 46-67.



- Hilpert, Martin. 2012. "Diachronic collocation analysis meets the noun phrase: Studying many a noun in COHA." In Terttu Nevalainen and Elizabeth Closs Traugott (eds.), *The Oxford handbook of the history of English*. Oxford: Oxford University Press. 233-44.
- Hofland, K. & S. Johansson. 1982. *Word frequencies in British and American English*. Bergen: Norwegian Computing Centre for the Humanities/London: Longman.
- Leech, Geoffrey & Roger Fallon. 1992. Computer corpora—What do they tell us about culture? *ICAME Journal* 16. 29-50.
- Michel, Jean-Baptiste, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew K. Gray, The Google Books Team, Joseph P. Pickett, Dale Hoiberg, Dan Clancy, Peter Norvig, Jon Orwant, Steven Pinker, Martin A. Nowak, and Erez Lieberman Aiden. (2011) Quantitative Analysis of Culture Using Millions of Digitized Books. *Science* 331: 176-182.
- Oakes, Michael & Malcolm Farrow. 2007. Use of the chi-square test to examine vocabulary differences in English-language corpora representing seven different countries. *Literary and Linguistic Computing* 22(1). 85-100.
- Sigley, Robert & Janet Holmes. 2002. Looking at *girls* in corpora of English. *Journal of English Linguistics* 30(2). 138-157.