





Introduction to Text Corpora and Their Applications

History and development of language corpora

Lucie Chlumská, Ph.D.

lucie.chlumska@korpus.cz





OUTLINE:

1. LECTURE

- where did the 'corpus' come from in the first place?
- pre-electronic corpora and the rise of computers
- 1st generation corpora
- 2nd generation corpora and beyond

2. SEMINAR

- reading (W. Teubert): the nature of language
- what is a language...? what is a discourse...?





LECTURE





The origins of corpus and corpus linguistics



Corpus is...

<http://www.oed.com/view/Entry/41873?redirectedFrom=corpus&>

Is corpus just a text archive?

Corpus is:

- assembled with particular purposes in mind and often *representative* of some language or text type (Leech 1992)
 - machine-readable and usually *annotated*
 - a *sample* of language



Corpus linguistics

Is corpus linguistics a discipline?

CL:

- does not concentrate on one single aspect of language (such as syntax, psycholinguistics or sociolinguistics etc.)
- focuses on **large amounts of authentic data** and strives to observe **patterns** in language (in the broadest sense of the word)
- works with **frequencies**, co-occurrences, collocation etc.
- focuses on **meaning** and context rather than grammar (unlike Chomsky)





Pre-electronic corpora



The beginnings

First attempts to collect data similar to corpora (before 1960s) were made in the following areas:

- biblical and literary studies
- lexicography
- dialect studies
- language education studies
- grammatical studies

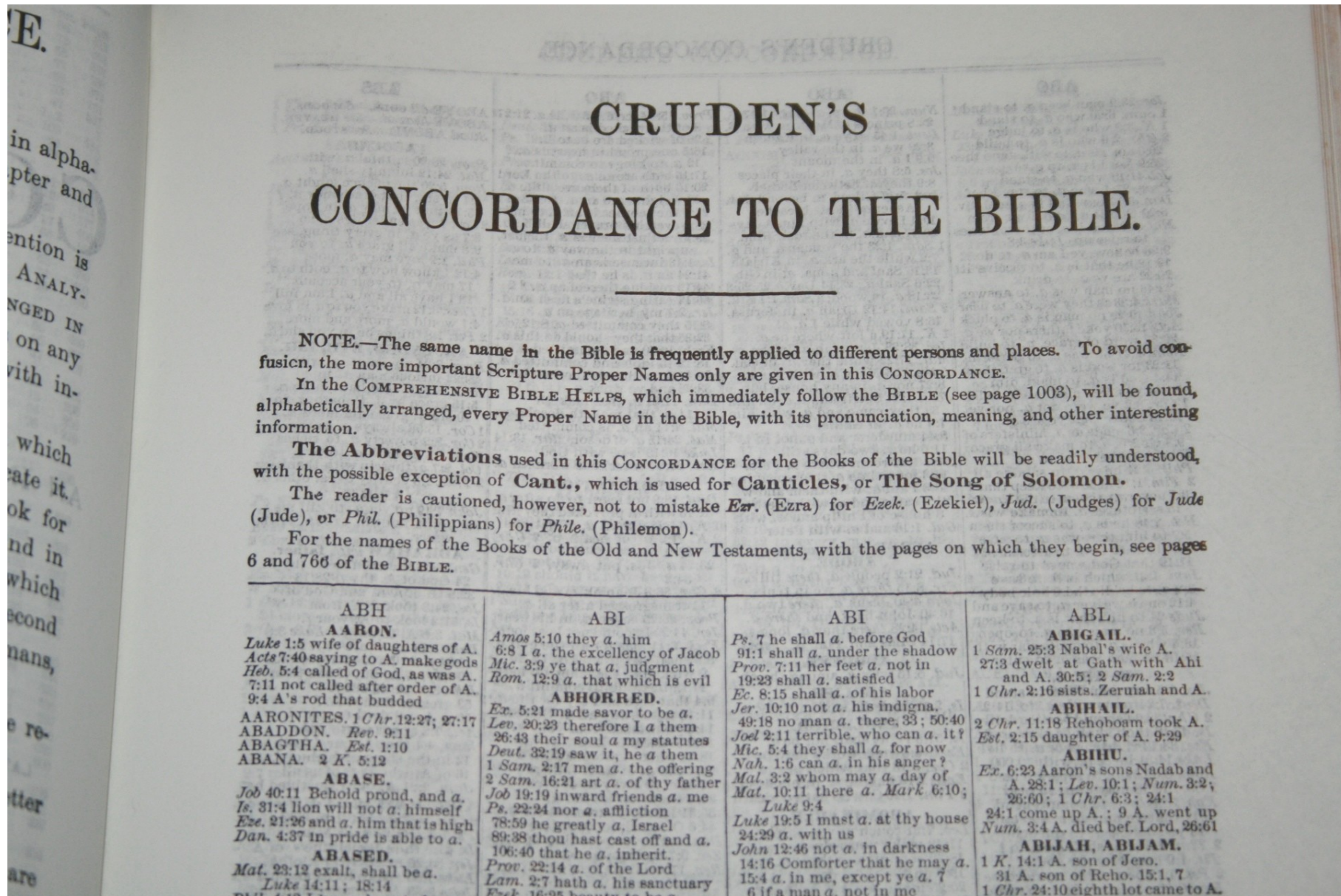


Biblical and literary studies

- Bible as a corpus on which to base commentaries or criticism, in the form of **word lists and concordances**
- in 18th century, **Alexander Cruden**, a London bookseller, proofreader and morals campaigner, produced the most famous of these for the Authorized (King James) Version of Bible
 - *Cruden's Concordance*, first published in 1736, 42 editions before 1879!
 - it included concordances not only for the main content words in the Bible, but also some function words and phrases (*how, you, he, once, how much, all the nations* etc.)



Biblical and literary studies



Lexicography

- as early as 17th century!
- **Samuel Johnson** recorded on slips of paper a large corpus of sentences from ‘writers of the first reputation’ to illustrate meanings and uses of English words in his *Dictionary of the English Language*
 - Johnson worked with 6 assistants to assemble over 150,000 illustrative citations for the app. 40,000 headword entries
- similarly, *Oxford English Dictionary (OED)* also corpus-based
 - twelfth and final volume published in 1928
 - 71 years of sustained work on a corpus of the canon of mainly literary written English from about AD 1000
 - 2,000 volunteer readers collected about five million citations amounting to 50 million words to illustrate 414,825 entries



Lexicography

- Compiling the *OED*...

- [How it began](#) 1857: The Philological Society of London calls for a new English Dictionary
- [More work than they thought](#) 1884: Five years into a proposed ten-year project, the editors reach *ant*
- [One step at a time](#) 1884-1928: The Dictionary is published in fascicles
- [Keeping it current](#) 1933-1986: Supplements to the *OED*
- [Making it modern](#) 1980s: The Supplements are integrated with the *OED* to produce its Second Edition
- [Into the electronic age](#) 1992: The first CD-ROM version of the *OED* is published
- [The future has begun](#) The present: The *OED* is now being fully revised, with new material published in parts online



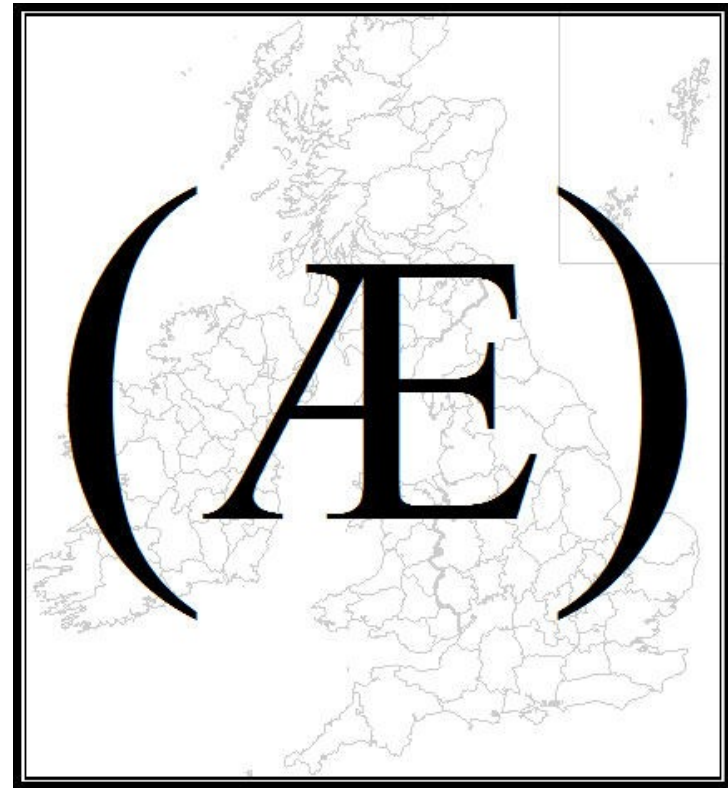
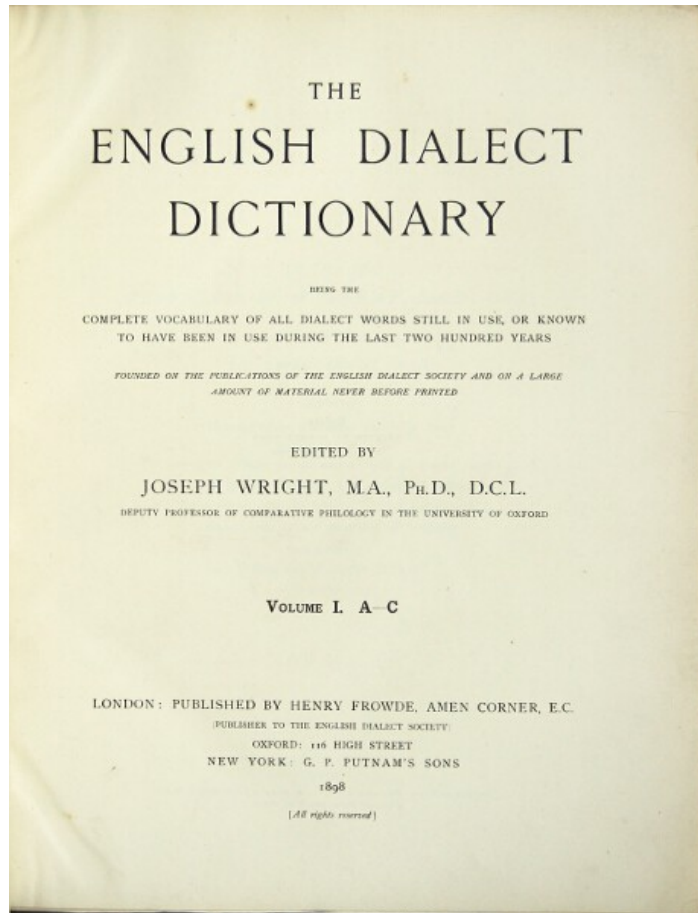
Lexicography

- parallel to the work on the second edition of *OED* in the latter part of 19th century, another great corpus of citations was being assembled to support the third edition of Noah Webster's *An American Dictionary of the English Language*
- in 1961, the third edition of Webster's *New International Dictionary* had available a corpus of over 10 million citation slips
- probably the last major English dictionary to be completed without an electronic database...



Dialect

- interest in linguistic variation in regional dialects in 19th century
 - the *English Dialect Dictionary* (Wright, 1898–1905)
 - *The Existing Phonology of English Dialects* (A. Ellis, 1889)



Language education

- some of the most influential corpus-based research in the 1st half of 20th century had a **pedagogical purpose**
 - **Edward Thorndike** (1921) compiled a corpus of 4.5 million words from 41 sources to make a frequency list (for better curricula materials for teaching literacy to native speakers of English in the US)
 - $\frac{3}{4}$ Bible and classic fiction, $\frac{1}{4}$ letters, newspapers, school readers
 - a corpus to gather statistical information of German words and letters in order to improve the stenographers' training
 - **J. W. Kaeding** needed help of 5,000 assistants over a period of years to process the corpus of 11 million words he used in 1890s for his analysis



Grammar

- **Otto Jespersen**, Danish professor, is said to have his villa filled with shoeboxes containing hundreds of thousands of paper slips with examples of interesting English sentences
 - monumental work *A Modern English Grammar on Historical Principles* (1909–49)
- **Charles C. Fries** used a corpus of letters written to the US government by persons of different educational and social backgrounds to demonstrate social class differences in usage in his *American English Grammar* (1940)
- Later in *The Structure of English* (1952) he used a 250,000-word corpus of recorded telephone conversations
- he analyzed all his corpora manually...



Grammar

- the most important pre-electronic corpus was the *Survey of English Usage Corpus (SEU)* by *Randolph Quirk et al.* in 1968
- it marked a transition between earlier non-computerized corpus-based description and the rise of corpus linguistics
- founded in 1959 by Quirk, the SEU aimed to collect 200 samples (each about 5,000 words) representative of both written and spoken language > corpus of 1 million words
 - SEU Corpus contains texts produced between 1953 and 1987, originally available in the form of paper slips filed at the University College London
 - there was a slip for every word in the corpus, containing 17 lines of text plus a mark-up (grammatical features, prosody...)
 - basis for the *A Comprehensive Grammar of the English Language* (1985)





The history of electronic corpora



Three stages of electronic corpora building

1. **1960–1980** : learning how to build and maintain corpora up to a million words; no material available in electronic form
2. **1980–2000**:
 - the 1980s: decade of the scanner
 - the 1990s, or the *First Serendipity*, text becomes available as the by-product of computer type-setting
3. the **new millennium**, or the *Second Serendipity*, texts that never existed as hard copies become available in unlimited quantities from the internet



The first electronic corpora

- the Brown corpus
 - compiled in the 1960s at Brown University by Nelson Francis and Henry Kučera
 - a million words of American English from documents published in 1961
 - 500 samples of about 2,000 words
 - a broad range of categories (informative prose: 374 samples, imaginative prose: 126 samples)



The first electronic corpora

- the LOB corpus
 - 1970–1978 at the Lancaster and Oslo University
 - a comparable counterpart to the American Brown corpus
 - a million words of British English from documents published in 1961
 - 500 samples of about 2,000 words
 - better technology and tools (sentence-initial markers etc.)
- FROWN and FLOB: Freiburg University corpora from 1992



The second generation

- the BNC
 - 1991–1994 (BNC XML Edition from 2007)
 - 100-million word collection of samples of written and spoken English
 - monolingual, synchronic, general and sample corpus
 - 4,054 texts
 - 90 % written part, 10 % spoken part
 - spoken language: demographic and context-governed



Modern corpora

- larger, better annotated & lemmatized etc.
- monitor corpora
 - **COCA** (1990–2015): 520 million tokens
 - **NOW** (2010–yesterday): 5.1 billion+
- web corpora
 - include texts from the internet: **WaC, Aranea, EnTenTen, SYN, Gigafida** etc.



Size of corpora and related resources

- 1st generation = small (*1-5 million*): **Brown** corpus and Brown family corpora (LOB, Frown, FLOB)
- 2nd generation = moderately-sized (*100 million*): **BNC, SYN2000**
- larger, more up-to-date corpora (*cca 450 million*): **COCA**
- extremely large text archives (billions): Google Books
- the Web as Corpus

But does size matter?



Frequency of different phenomena

Table 1.2 *Frequency of different phenomena in COCA, BNC, and Brown (numbers explained in detail in Sections 3.1–3.5)*

	COCA (450 m)	BNC (100 m)	Brown (1 m)
1 Lexical: individual	(See discussion in Section 3.1 above)		
2 Lexical: word lists	100,705	43,758	3,956
3 Morphology: substrings	<i>-ousness</i> 112 <i>-ism</i> 512	<i>-ousness</i> 25 <i>-ism</i> 278	<i>-ousness</i> 1 <i>-ism</i> 6
4 Morphology: compare	<i>prove</i> {n/d} 2,616 + 3,001 <i>sincere</i> 85 + 65	<i>prove</i> {n/d} 82 + 1,169 <i>sincere</i> 11 + 12	<i>prove</i> {n/d} 3 + 7 <i>sincere</i> 1 + 0
5 Syntax: high frequency	modals 5,794k perfects 1,837k <i>be</i> passives 2,900k	modals 1,421k perfects 446k <i>be</i> passives 890k	modals 14k perfects 4k <i>be</i> passives 10k
6 Syntax: low frequency	<i>love</i> 12,178 + 5,393 <i>hate</i> 3,968 + 1,773 <i>for</i> 931	<i>love</i> 1,192 + 351 <i>hate</i> 389 + 475 <i>for</i> 103	<i>love</i> 10 + 2 <i>hate</i> 8 + 2 <i>for</i> 0
7 Phraseology: words	<i>true feelings</i> 654 <i>naked eye</i> 175	<i>true feelings</i> 148 <i>naked eye</i> 53	<i>true feelings</i> 2 <i>naked eye</i> 0
8 Phraseology: constructions	<i>way</i> 251v : 15,868t <i>into</i> 275v : 2,160t	<i>way</i> 83v : 3,533t <i>into</i> 111v : 358t	<i>way</i> 15v : 44t <i>into</i> 6v : 6t
9 Semantics: collocates	<i>riddle</i> (n) 57 <i>nibble</i> (v) 96 <i>crumbled</i> (j) 33 <i>serenely</i> (r) 24	<i>riddle</i> (n) 0 <i>nibble</i> (v) 13 <i>crumbled</i> (j) 1 <i>serenely</i> (r) 4	<i>riddle</i> (n) 0 <i>nibble</i> (v) 0 <i>crumbled</i> (j) 0 <i>serenely</i> (r) 0
10 Semantics: prosody	<i>budge</i> (v) 1,427 <i>cause</i> (v) 1,344	<i>budge</i> (v) 164 <i>cause</i> (v) 358	<i>budge</i> (v) 3 <i>cause</i> (v) 0





Let's talk language!





SEMINAR



Reading

- common reading:

Teubert, W. (2004). Corpus linguistics: a different look at language. In M. A. K. Halliday, W. Teubert, C. Yallop and A. Čermáková, *Lexicology and Corpus Linguistics*, pp 96-112. NY: Continuum.

- additional reading recommendations:

<http://www.grsampson.net/ARsy.pdf>

<http://www.vox.com/the-big-idea/2016/9/14/12910180/noam-chomsky-tom-wolfe-linguist>

<https://aeon.co/essays/the-evidence-is-in-there-is-no-language-instinct>



Discussion

- What is a language?
- What is meaning? How can we discover a meaning of a word?
- What do we understand under the term *discourse*?
- Can we say that: what cannot be found in a corpus, does not exist in language?
- Why did Noam Chomsky object against corpora and corpus linguistics?
- What is and is not corpus linguistics good for?
- What are the main differences between a dictionary and a corpus?
- What was first, written or spoken language?
- What is a word?





<https://www.futurelearn.com/courses/corpus-linguistics?lr=4>

