# ВСТУП ДО ПРИКЛАДНОЇ ЛІНГВІСТИКИ: КОРПУСНА ЛІНГВІСТИКА І ВИВЧЕННЯ МОВ

# Повторення матеріалу лекції 3

#### КОРПУСНА ЛІНГВІСТИКА: ПОНЯТТЯ «КОРПУС»

- корпусна лінгвістика (corpus linguistics; corpus linguist)
- комп'ютерна лексикографія (computational lexicography)

a corpus is a large and principled collection of natural texts.

(Biber, Conrad, & Reppen, 1998, p. 12)

A corpus is a collection of pieces of language text in electronic form, selected according to external criteria to represent, as far as possible, a language or language variety as a source of data for linguistic research.

(Sinclair, 2005)

A corpus is a collection of (1) machine readable (2) authentic texts (including transcripts of spoken data) which is (3) sampled to be (4) representative of a particular language or language variety.

(McEnery, Xiao, & Tono, 2006, p. 5)

# ПРАКТИЧНЕ ЗАНЯТТЯ КОРПУСНА ЛІНГВІСТИКА І ВИВЧЕННЯ МОВ

ПРОБЛЕМИ створення електронних словників на основі друкованих продуктів:

- Ідеалізація слів і тенденція до науковості визначень (spider "member of Arachnidae, class Aranea)
- Ігнорування динаміки лексикону (неологізм *appendicitis* виключено зі словника Murray (1870)
- Невизначеність щодо власних назв (Shakespeare, England vs Persil, Pepsi)
- Складність диференціяції значень (*to pour*)
- Ілюзорність меж слів та складність їх розмужування (of course, insofar as, fire truck)



### Завдання 1:

# Реструктурація словників у корпусний формат

- Розгляньте корпуси та визначте:
- Який матеріал вони містять.
- Що слугувало джерелами матеріалу.
- Яким чином матеріал організований.
- Як можна застосувати корпусні дані для дослідження або вивчення мови.



#### **COCA Corpus of Contemporary American English**

https://www.english-corpora.org/coca/

#### The Student-Transcribed Corpus of Spoken American English

https://www.spokencorpus.org

Ta інші <a href="https://www.english-corpora.org">https://www.english-corpora.org</a>

Corpus		Download	# words	Dialect	Time period	Genre(s)
News on the Web (NOW)		0	23.1 billion+	20 countries	20 countries 2010- <b>yesterday</b>	
iWeb: The Intelligent Web-based	0	14 billion	6 countries	2017	Web	
Global Web-Based English (GloWbE)		0	1.9 billion	20 countries	2012-2013	Web (incl blogs)
Wikipedia Corpus		•	1.9 billion	(Various)	2014	Wikipedia
Coronavirus Corpus		0	1.5 <b>billion</b> 20 countrie		2020-2023	Web: News
Corpus of Contemporary American English (COCA)		0	1.0 billion	American	1990-2019	Balanced
The only large, recent, genre of English, as well as the mos		re-balanced corpus ost widely-used	475 million	American	American 1820-2019	
The TV Corpus	online corpus	325 million		6 countries	1950-2018	TV shows
The Movie Corpus		0	200 million 6 countr		1930-2018	Movies
Corpus of American Soap Operas		0	100 million American		2001-2012	TV shows
Hansard Corpus			1.6 <b>billion</b>	British	1803-2005	Parliament
Early English Books Online (EEBO)			755 million British		1470s-1690s	(Various)
Corpus of US Supreme Court Opinions			130 million	American	1790s-2017	Legal opinions
TIME Magazine Corpus			100 million	American	1923-2006	Magazine
British National Corpus (BNC) *			100 million	British	1980s-1993	Balanced
Strathy Corpus (Canada)			50 million	Canadian	1970s-2000s	Balanced
CORE Corpus			50 million	6 countries	2014	Web
From Google Books n-grams (compare)						
American English			155 billion	American	1500s-2000s	(Various)
British English			34 billion	British	1500s-2000	(Various)



#### **Corpus of Contemporary American English**















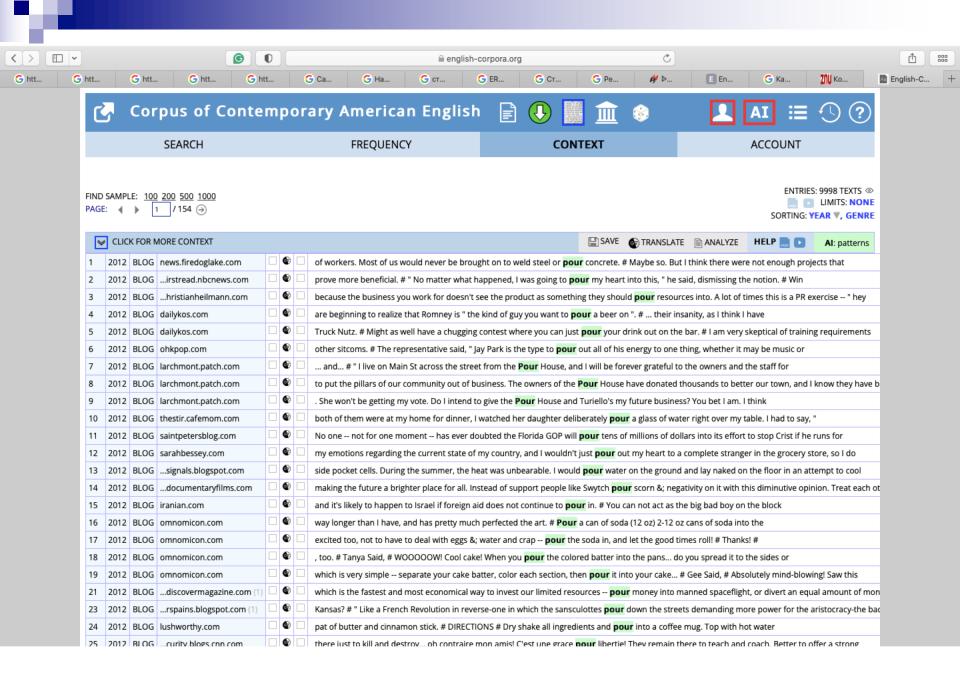


SEARCH FREQUENCY CONTEXT TEXTS

Download spreadsheet for all 485,179 texts (with summary by year, genre, and sub-genre)

The corpus is composed of more than **one billion words** in 485,202 texts, including 24-25 million words each year from 1990-2019. For each year (and therefore overall, as well), the corpus is evenly divided between the genres of TV and Movies subtitles, spoken, fiction, popular magazines, newspapers, and academic journals. This is important, because if you want to compare different years, you need to be comparing "apples" to "apples" (i.e. same genre balance in the different periods).

YEAR	BLOG	WEB	TV / MOVIES	SPOKEN	FICTION	MAGAZINE	NEWSPAPER	ACADEMIC	TOTAL
TOTAL	125,496,215	129,899,426	128,013,334	127,396,916	119,505,292	127,352,014	122,959,393	120,988,348	1,001,610,938
1990			3,207,900	4,374,469	4,162,242	4,101,447	4,082,931	3,983,143	23,912,132
1991			3,379,151	4,316,898	4,192,646	4,209,838	4,104,806	4,051,046	24,254,385
1992			3,183,858	4,523,054	3,893,956	4,288,694	4,092,031	4,028,147	24,009,740
1993			3,785,924	4,487,978	3,921,244	4,254,351	4,153,070	4,150,671	24,753,238
1994			4,375,338	4,457,726	3,870,757	4,310,375	4,147,947	4,047,115	25,209,258
1995			5,006,966	4,548,602	3,846,412	4,314,737	4,122,703	4,016,371	25,855,791
1996			4,384,976	4,095,266	3,758,787	4,338,766	4,099,305	4,110,209	24,787,309
1997			4,380,670	3,904,996	3,617,741	4,368,917	4,153,906	4,420,786	24,847,016
1998			4,390,197	4,446,217	3,779,801	4,393,835	4,122,295	4,111,453	25,243,798
1999			4,381,144	4,445,564	4,154,537	4,391,146	4,107,423	4,023,282	25,503,096



#### FINDING KEYWORDS FOR A WORD OR PHRASE WITH KWIC ENTRIES

The corpora can help you to quickly and easily find words related to a given word or phrase. For example, for single words (such as enzyme, lighthouse, environmental, or sew) you can find collocates (nearby words) and topics (words that co-occur anywhere in the text). And using Virtual Corpora, you can quickly and easily create a "sub-corpus" of texts that contain a given word or phrase (such as New York, Harry Potter, investment, or refugee).

Another easy way to find words related to a particular word or phrase is via the "Analyze (Text)" function from the Keyword in Context page (KWIC; concordances). The following are the steps to do this.

Do a LIST search for a word or phrase, and then click on SAMPLE in the LIST display.



This will give you randomly-selected KWIC lines. Click on Analyze [1], which will show the Analyze function [2]. For the best results, click the checkbox [3] to select all of the KWIC lines.



This then creates a "mini-text" of all of the KWIC lines. If there are about 30 words per line and 100 lines of text, this is a text of about 3000 words, and you can analyze this as you would any other text. You can click on any word in the "text" [1] or click on one of the words in the keyword list [2].



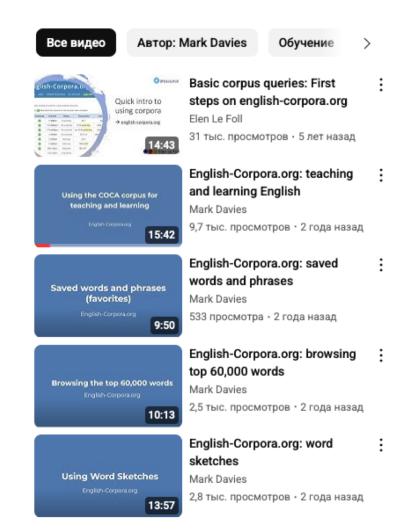
For any word that you click on (such as *habitat*, above), you can see a detailed <u>word sketch</u>, which shows (among other things) frequency information [1], definitions [2], links to external pronunciation, videos, and images [3], synonyms and semantically-related words [4], related topics [5], collocates [6], and Keyword in Context entries (not shown here)



And all of this can be done is just 5-10 seconds, via the "Analyze" function in the Keyword in Context (KWIC) display.

# Корисні відео: як працювати з СОСА

https://www.youtube.com/watch?v=\_\_G2PG46180&t=1s



# Корисні відео: як працювати з СОСА

https://www.youtube.com/watch?v=OmDsGCJZMx8



English-Corpora.org: teaching and learning English











