# COCA

*A Practical and Very Much Helpful Web-based Application for Academic Writing*

**Aysel Saricaoglu**
**English Language Education**

TED ÜNİVERSİTESİ

# What does "corpus/corpora" mean?

Blog    BuzzWord    Open Dictionary    Games    Resources    API    More

MACMILLAN DICTIONARY

Search

**Did you know?**

Click any word in a definition or example to find the entry for that word

## corpus - definition and synonyms

Show less

NOUN [COUNTABLE]    Pronunciation    /ˈkɔː(r)pəs/    Word Forms

*Using the thesaurus*

Contribute to our Open Dictionary

1 FORMAL a collection of writing, for example all the writings of one person

Synonyms and related words

**Collections, stores and sets of things:**
*agglomeration, arsenal, assemblage...*

Explore Thesaurus

**Related words**

habeas corpus NOUN

2 LINGUISTICS a collection of written and spoken language stored on computer and used for language research and writing dictionaries

Synonyms and related words

**Dictionaries and making dictionaries:**
*concordance, corpus, definition...*

In Latin: "body" (same root as "corpse")

Linguistically: a large (sampled to be representative) and usually electronic (machine-/computer-readable) collection of authentic texts (naturally occurring / written by native speakers)

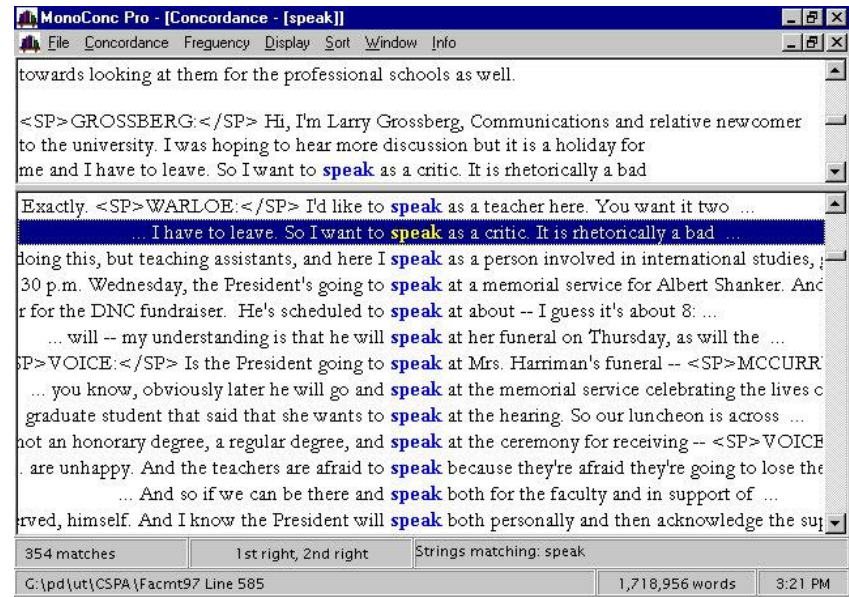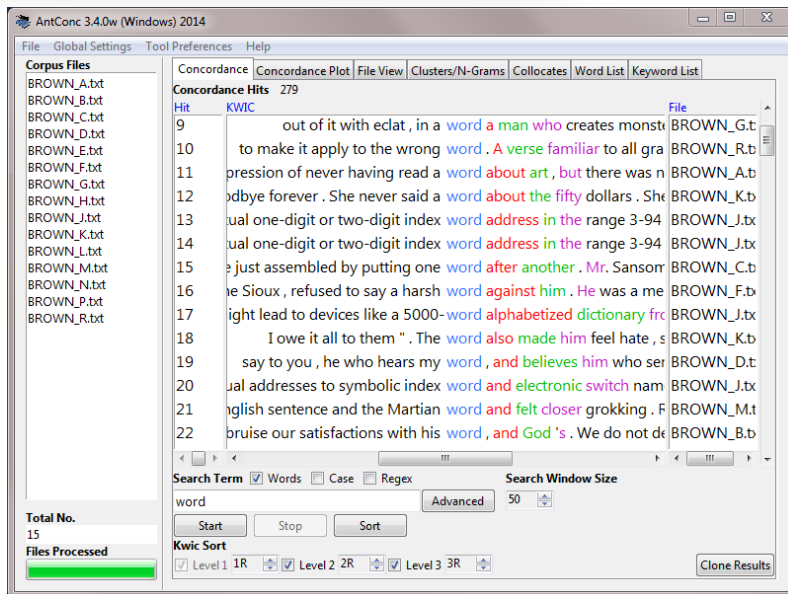Self-compilable

E.g.,

     - a corpus of 90s songs

     - a corpus of civil engineering research articles

     - a corpus of newspapers

     - a corpus of telephone conversations

     - a corpus of air traffic control conversations

Corpus linguistics: the use of corpora for language study

# You don't have to be a computational linguist or a computer engineer to use corpora: you can use *concordancers*

# Why use corpora?

- to find out and model how native speakers speak and write

- to carry out linguistic analyses

- to analyze variation and change in different languages

- to develop authentic language teaching materials

- to  get precise data on the target languages for translation purposes

- to look at changes in culture and society

- to use data in natural language processing projects

5. Children who are securely attached to their mother by the time they enter their second year of life are better equipped to ---- new experiences and relationships.

A) depend on
B) turn down

C) refrain from
D) cope with

E) carry out

DOĞRU YANIT: D

2) Managing traffic flow at peak periods and dealing with incidents, such as crashes, are ---- problems for transport planners.

A) challenging

B) instructive

C) accessible

D) favorable

E) functional

DOĞRU CEVAP: A

# COCA: The Corpus of Contemporary American English

- created by Lingustics professor Mark Davies, BYU

- freely-available

- more than 520 million words of texts between 1990-2015

- equally divided among spoken, fiction, popular magazines, newspapers, and academic texts

# Let's try it out!

## http://corpus.byu.edu/coca/

Matching forms + frequencies

– Search for **how individual words are used**: innovative

– Search for **how phrases are used**: freak out, FREAK OUT

– Search for **all forms of a word**: EAT; UNDERSTAND

– Search with **parts of speech** : well-established NOUN (or well-established [nn*]); speak ADV; speak [rp*]; poached NOUN

– Search for **patterns matching or including a word (wildcards)**:  *new*

– Search for **any word in phrases**:  as * to; as * by

– Search by **synonyms**: =fundamentally; in =particular NOUN

– Search with **alternants**: ADJ perspectives|outcomes; essential|necessary NOUN

– Search by **excluding**: pretty –NOUN

– Search by **combining**: =significant findings; we =aim to; these findings =show

Section-/Genre-  & time-specific matching forms + frequencies

– Search for **lexical variation between genres:** You have to know; sort of; type of; the cream of the crop; the best ; busted; unattractive ; lots of; [get] [vvn*]; look up; research

– Search for **lexical change over time:** BE likely a|the; Recep Tayyip Erdogan; Ataturk; Ankara; Middle East Technical University

# Collocates function
*words that occur with other words; nearby words*

- Search by **distance and grammatical category**: expand (word) NOUN (collocates) | mistakenly (word) VERB (collocates) | opportunity (word) ADJ (collocates)

- Search **by distance only for any word**: controversial (word) * (collocates)

- Search **by distance and for synonyms of a word**: opportunity (word) =great (collocates)

– Search for **differences in usage and meaning**: big/further /less (word1) tall/farther /fewer

(word2) NOUN (collocates) |man (word1) woman (word2) ADJ (collocates)

– Search by **sorting the words to the left or right**: LOOK into (123R) | money (321) | sick of (123R)

# Thank you!

*Questions/Comments?*

Aysel Saricaoglu
aysel.saricaoglu@tedu.edu.tr

www.tedu.edu.tr