

# Format download file *Leipzig Corpora Collection*

This file describes the format of the download corpus files of the Leipzig Corpora Collection. All files are encoded in UTF-8. Columns are separated by tabs.

## Word list

*The file contains the word list of all word forms of the corpus. Words are ordered by their frequency in descending order. The first 100 IDs of the word list are reserved for special characters.*

Filename: \*\_words.txt

Format: Word\_ID Word Word Frequency

## Sentences list

*The file contains all sentences of the corpus.*

Filename: \*\_sentences.txt

Format: Sentence\_ID Sentence

## Sources list

*The file contains information about the used sources.*

Filename: \*\_sources.txt

Format: Source\_ID Source Date

## Neighbourhood cooccurrences

*The file contains information about how often two words occurred in direct neighbourhood in the the corpus and the significance of those cooccurrences based on log-likelihood. In the file, word1 occurs immediately left of word2.*

Filename: \*\_co\_n.txt

Format: Word1\_ID Word2\_ID Number\_of\_Cooccurrences Significance

## Sentences cooccurrences

*The file contains information about how often two words occurred in the same sentence and the significance of those cooccurrences based on log-likelihood.*

Filename: \*\_co\_s.txt

Format: Word1\_ID Word2\_ID Number\_of\_Cooccurrences Significance

## Inverted list

*The file contains information about the occurrences of words in sentences (and optional their position in the sentence).*

Filename: \*\_inv\_w.txt

Format: Word\_ID Sentence\_ID (Position\_in\_Sentence)

## Inverted source list

*The file contains the mapping of a sentence to the sources from which it was extracted.*

Filename: \*\_inv\_so.txt

Format: Source\_ID Sentence\_ID

## Metadata

*The file contains several metadata about the creation process of the corpus.*

Filename: \*\_meta.txt

Format: Metadaten\_ID Key Value

## Import script

*The import script can be used to import the files into a MySQL database.*

Filename: \*-import.sql

Example (Linux): `$ mysql Database_Name < Database_Name-import.sql`