

Format Download „Deutscher Wortschatz“

Im Folgenden werden alle Angaben kurz erläutert, die Bestandteil einer Downloaddatei des Wortschatz-Projektes sind. Alle Dateien sind in UTF-8 kodiert. Spalten werden durch Tabulatoren getrennt.

Wortliste

Die Datei enthält die Wortliste aller im Korpus vorkommenden Wörter. Wörter sind absteigend sortiert nach ihrer Frequenz. Die ersten 100 Wort-IDs sind für Sonderzeichen reserviert.

Dateiname: *_words.txt

Format: Wort_ID Wort Wort Frequenz

Satzliste

Die Datei enthält alle Sätze des Korpus.

Dateiname: *_sentences.txt

Format: Satz_ID Satz

Quellenliste

Die Datei enthält Angaben zu den verwendeten Quellen.

Dateiname: *_sources.txt

Format: Quellen_ID Quelle Datum

Nachbarschaftskookkurrenzen

Die Datei enthält Informationen wie oft zwei Wörter in unmittelbarer Nachbarschaft im Korpus vorkommen und die Signifikanz dieses Auftretens auf Basis von Log-Likelihood. Wort1 steht dabei unmittelbar links von Wort2.

Dateiname: *_co_n.txt

Format: Wort1_ID Wort2_ID Anzahl_Vorkommen Signifikanz

Satzkookkurrenzen

Die Datei enthält Informationen wie oft zwei Wörter im gleichen Satz im Korpus vorkommen und die Signifikanz dieses Auftretens auf Basis von Log-Likelihood.

Dateiname: *_co_s.txt

Format: Wort1_ID Wort2_ID Anzahl_Vorkommen Signifikanz

Inverse Liste

Die Datei enthält die Zuordnung eines Wortes zu den Sätzen in welchem es (und optional an welcher Position in diesem) vorkommt.

Dateiname: *_inv_w.txt

Format: Wort_ID Satz_ID (Position_im_Satz)

Inverse Quellenliste

Die Datei enthält die Zuordnung eines Satzes zur Quelle aus der er extrahiert wurde.

Dateiname: *_inv_so.txt

Format: Satz_ID Quellen_ID

Allgemeine Metadaten

Die Datei enthält verschiedene Metadaten zum Erstellungsprozess des Korpus.

Dateiname: *_meta.txt

Format: Metadaten_ID Key Value

Importskript

Das Importskript kann zum Import der Dateien in eine MySQL-Datenbank genutzt werden.

Dateiname: *-import.sql

Beispielaufruf (Linux): `$ mysql Datenbank_Name < Datenbank_Name-import.sql`