DH 1: Sprache und Text

# Korpusindexierung und -abfragen

**Andreas Blombach, Stephanie Evert**
Lehrstuhl für Korpus- und
Computerlinguistik
https://www.linguistik.phil.fau.de

CL FAU **Friedrich-Alexander-Universität**
**Philosophische Fakultät und**
**Fachbereich Theologie**

# Repräsentationsformate & Indexierung

# Korpusannotation: Rohtext + Metadaten

It seemed a day much as any other until I happened to look out of the back window.  There was a little garden behind the house; a well-mown lawn surrounded by a neatly cut hedge, a few bushes and colourful flowers.

**Metadaten**
title:          The Garden
author:         Stefan Evert
author sex:     male
date:           05.08.1991

# Korpusannotation: Tokenisierung

It seemed a day much as any other until I happened to look out of the back window .  There was a little garden behind the house ; a well-mown lawn surrounded by a neatly cut hedge , a few bushes and colourful flowers .

# Korpusannotation: Satzsegmentierung

\<s\> It seemed a day much as any other until I happened to look out of the back window . \</s\>

\<s\> There was a little garden behind the house ; a well-mown lawn surrounded by a neatly cut hedge , a few bushes and colourful flowers . \</s\>

# Korpusannotation: POS-Tagging (Wortartenannotation)

<s> It$_{PP}$ seemed$_{VBD}$ a$_{DT}$ day$_{NN}$ much$_{RB}$ as$_{IN}$ any$_{DT}$ other$_{JJ}$ until$_{IN}$ I$_{PP}$ happened$_{VBD}$ to$_{TO}$ look$_{VB}$ out$_{RP}$ of$_{IN}$ the$_{DT}$ back$_{JJ}$ window$_{NN}$ .$_{SENT}$ </s>

<s> There$_{EX}$ was$_{VBD}$ a$_{DT}$ little$_{JJ}$ garden$_{NN}$ behind$_{IN}$ the$_{DT}$ house$_{NN}$ ;$_{:}$ a$_{DT}$ well-mown$_{VBN}$ lawn$_{NN}$ surrounded$_{VBN}$ by$_{IN}$ a$_{DT}$ neatly$_{RB}$ cut$_{VBN}$ hedge$_{NN}$ ,$_{,}$ a$_{DT}$ few$_{JJ}$ bushes$_{NNS}$ and$_{CC}$ colourful$_{JJ}$ flowers$_{NNS}$ .$_{SENT}$ </s>

# Korpusannotation: Lemmatisierung

\<s> It_PP ^it seemed_VBD ^seem a_DT ^a day_NN ^day much_RB ^much as_IN ^as any_DT ^any other_JJ ^other until_IN ^until I_PP ^I happened_VBD ^happen to_TO ^to look_VB ^look out_RP ^out of_IN ^of the_DT ^the back_JJ ^back window_NN ^window ._SENT ^. \</s>

\<s> There_EX ^there was_VBD ^be a_DT ^a little_JJ ^little garden_NN ^garden behind_IN ^behind the_DT ^the house_NN ^house ;_: ^; a_DT ^a well-mown_VBN ^well-mown lawn_NN ^lawn surrounded_VBN ^surround by_IN ^by a_DT ^a neatly_RB ^neatly cut_VBN ^cut hedge_NN ^hedge ,_, ^, a_DT ^a few_JJ ^few bushes_NNS ^bush and_CC ^and colourful_JJ ^colorful flowers_NNS ^flower ._SENT ^. \</s>

# XML-Markup der Annotation

## Standard für Datenaustausch und -archivierung

```xml
<?xml version="1.0" encoding="UTF-8"?>          ⟵  XML-Deklaration
<corpus>
  <story title="The Garden">
    <p>
      <s>
        <token pos="PP"  lemma="it">It</token>
        <token pos="VBD" lemma="seem">seemed</token>
        <token pos="DT"  lemma="a">a</token>
        <token pos="NN"  lemma="day">day</token>
        <token pos="RB"  lemma="much">much</token>
        <token pos="IN"  lemma="as">as</token>
        <token pos="DT"  lemma="any">any</token>
        <token pos="JJ"  lemma="other">other</token>
        <token pos="IN"  lemma="until">until</token>
        <token pos="PP"  lemma="I">I</token>
        ...
      </s>
    </p>
  </story>
</corpus>
```

# XML-Markup der Annotation

```xml
<?xml version="1.0" encoding="UTF-8"?>
<corpus>                    ← Wurzelelement / root element
  <story title="The Garden">
    <p>                     ← Start-Tag des XML-Elements
      <s>
        <token pos="PP"  lemma="it">It</token>    ← Element mit Attributen
        <token pos="VBD" lemma="seem">seemed</token>
        <token pos="DT"  lemma="a">a</token>
        <token pos="NN"  lemma="day">day</token>
        <token pos="RB"  lemma="much">much</token>
        <token pos="IN"  lemma="as">as</token>
        <token pos="DT"  lemma="any">any</token>
        <token pos="JJ"  lemma="other">other</token>
        <token pos="IN"  lemma="until">until</token>
        <token pos="PP"  lemma="I">I</token>
        ...
      </s>                  ← korrespondierendes End-Tag
    </p>
  </story>
</corpus>
```

# XML-Dokument = geordneter Baum

# XML als Repräsentationsformat



```
H9C.xml*  ×

1  ▽ <bncDoc, xml:id="H9C">
2  ▽     <teiHeader>                                    ← TEI-Header enthält Metadaten
3  ▽         <fileDesc>
4  ▽             <titleStmt>
5  ▽                 <title> The prince of darkness. Sample containing about 44223 words from a book
6                         (domain: imaginative) </title>
7  ▽                 <respStmt>
8                         <resp> Data capture and transcription </resp>
9                         <name> Oxford University Press </name>
10                    </respStmt>
11                </titleStmt>
12 ▽             <editionStmt>                           Beispiel aus dem British National Corpus
13                    <edition>BNC XML Edition, December 2006</edition>
14                </editionStmt>
15                <extent> 44223 tokens; 44797 w-units; 3933 s-units </extent>
16 ▽             <publicationStmt>
17 ▽                 <distributor>Distributed under licence by Oxford University Computing Services on
18                        behalf of the BNC Consortium.</distributor>
19 ▽                 <availability> This material is protected by international copyright laws and may
20                        not be copied or redistributed in any way. Consult the BNC Web Site at
21                        http://www.natcorp.ox.ac.uk for full licencing and distribution
22                        conditions.</availability>
23                    <idno type="bnc">H9C</idno>
24                    <idno type="old"> PDarkn </idno>
25                </publicationStmt>
26 ▽             <sourceDesc>
27 ▽                 <bibl>
28                        <title>The prince of darkness. </title>
29                        <author domicile="Epping" n="DoherP1">Doherty, P C</author>
30 ▽                     <imprint n="HEADLI1">
31                            <publisher>Headline Book Publishing plc</publisher>
32                            <pubPlace>London</pubPlace>
33                            <date value="1992">1992</date>
34                        </imprint>
35                    </bibl>
36                </sourceDesc>
37            </fileDesc>
38 ▽         <encodingDesc>
39 ▽             <tagsDecl>
40 ▽                 <namespace name="">
```

Informationen über den Text

11

# XML als Repräsentationsformat

TEI-Body enthält Objektdaten

Textstruktur & Darstellung

Token mit Annotationen

```
80  <wtext type="FICTION">
81      <pb n="69"/>
82      <div level="1">
83          <head>
84              <s n="2">
85                  <w c5="NN1" hw="chapter" pos="SUBST">Chapter </w>
86                  <w c5="CRD" hw="5" pos="ADJ">5</w>
87              </s>
88          </head>
89          <p>
90              <s n="3">
91                  <w c5="VVB-NN1" hw="ranulf" pos="VERB">Ranulf </w>
92                  <w c5="CJC" hw="and" pos="CONJ">and </w>
93                  <w c5="NP0" hw="dame" pos="SUBST">Dame </w>
94                  <w c5="NP0" hw="agatha" pos="SUBST">Agatha </w>
95                  <w c5="VBD" hw="be" pos="VERB">were </w>
96                  <w c5="VVG" hw="wait" pos="VERB">waiting </w>
97                  <w c5="PRP" hw="for" pos="PREP">for </w>
98                  <w c5="PNP" hw="he" pos="PRON">him </w>
99                  <w c5="PRP" hw="near" pos="PREP">near </w>
100                 <w c5="AT0" hw="the" pos="ART">the </w>
101                 <w c5="NN1-NP0" hw="galilee" pos="SUBST">Galilee </w>
102                 <w c5="NN1" hw="gate" pos="SUBST">Gate</w>
103                 <c c5="PUN">, </c>
104                 <w c5="AT0" hw="the" pos="ART">the </w>
105                 <w c5="AJ0" hw="young" pos="ADJ">young </w>
106                 <w c5="NN1" hw="nun" pos="SUBST">nun </w>
107                 <w c5="AV0" hw="apparently" pos="ADV">apparently </w>
108                 <w c5="VVG" hw="enjoy" pos="VERB">enjoying </w>
109                 <w c5="AT0" hw="an" pos="ART">an </w>
110                 <w c5="NN1" hw="account" pos="SUBST">account </w>
111                 <w c5="PRF" hw="of" pos="PREP">of </w>
112                 <w c5="CRD" hw="one" pos="ADJ">one </w>
113                 <w c5="PRF" hw="of" pos="PREP">of </w>
114                 <w c5="DPS" hw="he" pos="PRON">his </w>
115                 <w c5="NN1" hw="manservant" pos="SUBST">manservant</w>
116                 <w c5="POS" hw="'s" pos="UNC">'s </w>
117                 <w c5="DT0" hw="many" pos="ADJ">many </w>
118                 <w c5="NN2" hw="escapade" pos="SUBST">escapades </w>
119                 <w c5="PRP" hw="in" pos="PREP">in </w>
120                 <w c5="NP0" hw="london" pos="SUBST">London</w>
121                 <c c5="PUN">.</c>
```

**XML-Prinzip:**
Durch Entfernen aller XML-Tags kann der ursprüngliche Objekttext wiederhergestellt werden

# XML-Standards

- XML (*Extensible Markup Language*) ist ein weitverbreiteter Standard für hierarchisch strukturierte Annotationen

- ein wohlgeformtes XML-Document legt nur die Struktur der Annotation/Auszeichnung fest, nicht die Semantik (also was was bedeutet)

- eine DTD (*document type declaration*) oder ein XML-Schema legt gültige Element- und Attributnamen fest
  - … erklärt aber ohne Dokumentation noch immer nicht die Semantik!

- Austauschformate für Textkorpora:
  TEI (*Text Encoding Initiative*), XCES (*Corpus Encoding Standard*),
  ISO 24612: LAF (Linguistic Annotation Framework)
  - ideal für Archivierung und Interoperabilität
  - für Korpusabfragen u.a. wird aber effizientere Implementierung benötigt

# CWB

- IMS Open Corpus Workbench
    - ursprünglich 1993–1996 entwickelt (IMS Stuttgart)
    - Anwendungen: statistische Sprachverarbeitung, Lexikographie, Korpuslinguistik
    - Open-Source-Veröffentlichung 2005 (GPL)
    - Aktuelle Version: CWB 3.5 (UTF-8, Korpora bis zu 2 Milliarden Wörter)
- Standardoberfläche: CQPweb
    - diverse simplere und/oder spezialisiertere Web-Oberflächen verfügbar
    - Kommandozeilen-CQP für erfahrene Benutzer (Uli Heid)
- SketchEngine: gleiches Datenmodell, gleiche Abfragesyntax
    - aber unterschiedliche Implementierung

http://cwb.sf.net/

# Vertical text format (.vrt)

**Einfacheres, effizienteres Format → wird von CWB & NLP-Tools verwendet**

```
<corpus>
<story title="The Garden">
<p>
<s>
It        PP      it
seemed    VBD     seem
a         DT      a
day       NN      day
much      RB      much
as        IN      as
any       DT      any
other     JJ      other
until     IN      until
I         PP      I
...
</s>
</p>
</story>
</corpus>
```

**metadata**
title:          The Garden
author:         Stefan Evert
author sex:     male
date:           05.08.1991

Tabulatorzeichen (\t, \x09)

# Vertical text format (.vrt)

## Textmetadaten in den XML-Start-Tags (nicht im Header!)

```
<corpus>
<text_title="The Garden" author="Stefan Evert" author_sex="male"
     date="1991-08-05">
<p num="1">
<s>
It        PP    it
seemed    VBD   seem
a         DT    a
day       NN    day
much      RB    much
as        IN    as
any       DT    any
other     JJ    other
until     IN    until
I         PP    I
...
</s>
</p>
</text>
</corpus>
```

CQPweb requires **<text>**,
SketchEngine prefers **<doc>**

sub-text level metadata

# CoNLL-Format(e)
## Vertical Text ohne Metadaten (weit verbreitet in NLP)

```
# story: "The Garden"
# paragraph #1
1    It        PP     it
2    seemed    VBD    seem
3    a         DT     a
4    fine      JJ     fine
5    day       NN     day
6    .         SENT   .

1    There     EX     there
2    was       VBD    be
3    an        DT     a
4    elephant  NN     elephant
5    .         SENT   .

# this is the end of the file
```

Metainformation in Kommentaren

Leerzeile = Satzgrenze

Token durchnummeriert (je Satz)

# CWB: Datenmodell

| #  | word     | pos | lemma   |
|----|----------|-----|---------|
| 0  | A        | DET | a       |
| 1  | fine     | ADJ | fine    |
| 2  | example  | NN  | example |
| 3  | .        | PUN | .       |
|    |          |     |         |
| 4  | Very     | ADV | very    |
| 5  | fine     | ADJ | fine    |
| 6  | examples | NN  | example |
| 7  | !        | PUN | !       |

implizite Nummerierung der Token:
Korpusposition (**cpos**)

# CWB: Datenmodell

| # | word | pos | lemma |
|---|------|-----|-------|
| (0) | `<text id="42" lang="English">` | | |
| (0) | `<s>` | | |
| 0 | A | DET | a |
| 1 | fine | ADJ | fine |
| 2 | example | NN | example |
| 3 | . | PUN | . |
| (3) | `</s>` | | |
| (4) | `<s>` | | |
| 4 | Very | ADV | very |
| 5 | fine | ADJ | fine |
| 6 | examples | NN | example |
| 7 | ! | PUN | ! |
| (7) | `</s>` | | |
| (7) | `</text>` | | |

XML-Tags als „unsichtbare" Token

# CWB: Datenmodell

| # | word | pos | lemma | | |
|---|---|---|---|---|---|

**p-Attribute**

| # | word | pos | lemma |
|---|---|---|---|
| (0) | `<text id="42" lang="English">` | | |
| (0) | `<s>` | | |
| 0 | A | DET | a |
| 1 | fine | ADJ | fine |
| 2 | example | NN | example |
| 3 | . | PUN | . |
| (3) | `</s>` | | |
| (4) | `<s>` | | |
| 4 | Very | ADV | very |
| 5 | fine | ADJ | fine |
| 6 | examples | NN | example |
| 7 | ! | PUN | ! |
| (7) | `</s>` | | |
| (7) | `</text>` | | |

id="42" lang="English"

**s-Attribute**

interne Dar-
stellung als
Tokenbereiche
(**cpos spans**)
der Form
[*start, end*]

# Korpusabfragen mit CQP

# Zum Ausprobieren: CQPweb & BNCweb

- **CQPweb** (flexible Web-Oberfläche für CWB)
  https://corpora.linguistik.uni-erlangen.de/cqpweb/
  - Login: `studentX` (`1` … `15`)
  - Passwort: `erlangen`
  - Dokumentation: Tutorial-Videos auf YouTube
    https://www.youtube.com/user/CorpusWorkbench

- **BNCweb** (speziell für British National Corpus)
  https://corpora.linguistik.uni-erlangen.de/bncweb/
  - gleicher Gäste-Login wie für CQPweb

# Simple query syntax

- viele Web-Oberflächen unterstützen eine „einfache" Abfragesyntax
  - Wort oder Wortfolge einfach direkt eingeben
  - Platzhalter (*wildcards*) für variable/optionale Elemente (eingeschränkt)
- hier: CEQL-Syntax (*Common Elementary Query Language*)
  - relativ mächtige einfache Abfragesprache
  - unterstützt von BNCweb, CQPweb und einigen anderen Web-Oberflächen
- Tutorial & Dokumentation
  - **Kap. 6** aus Hoffmann, Sebastian *et al.* (2008): *Corpus Linguistics with BNCweb – a Practical Guide*. Frankfurt a.M. [etc.]: Peter Lang.
  - CQPweb simple query manual
    https://cqpweb.lancs.ac.uk/doc/cqpweb-simple-syntax-help.pdf

# CEQL in a nutshell

- `speak`                        exakt diese Wortform
- `{speak}`                      alle flektierten Formen (Lemma)
- `at the end of the day`   exakt diese Wortfolge
- `is n't it \?`                 Tokenisierungsregeln & Escapes
- `*able`                        Suffix *-able*
- `+able`                        ohne das Wort *able* selbst
- `light_JJ`                     *light* als Adjektiv
- `Mr _N*`                       männl. Person (*Mr* + Substantiv)
- `[Mr,Mrs,Ms] _N*`              Person (männlich oder weiblich)
- `Mr _N* {be} _J*`             Aussage über die Person
  (wird später noch verbessert)

## CEQL in a nutshell

- `Smith:C`                          Groß- und Kleinschr. beachten
- `deja:d vu:d`                      Diakritika ignorieren
- `\D`                               Zahl (mind. eine Ziffer)
- `\u\u\u\u:C`                       Akronym (4 Großbuchstaben)
- `\u\L:C`                           nur erster Buchstabe groß
- `take * off`                       optionales Token
- `take ++*** off`                   zwei bis fünf optionale Token
- `in ( _JJ* )? time`               optionales Adjektiv
- `Mr (_N*)+ {be} (_RB)? _J*`        + = 1 oder mehr
- `his ( _JJS | most _JJ )* _N*`    Alternativen
- `<s> but`                          Satzanfang
- `<ne_type=PERSON> (+)+ </ne_type>` XML-Element

# CQP Query Syntax

- formale Abfragenotation
    - basiert auf regulären Ausdrücken auf mehreren Ebenen
    - ermöglicht es, das Suchmuster sehr präzise festzulegen
    - deutlich flexibler und mächtiger als die CEQL-Syntax
- von allen Web-Oberflächen unterstützt, die auf CWB basieren!
- Tutorial & Dokumentation
    - **Kap. 12** aus Hoffmann, Sebastian *et al.* (2008): *Corpus Linguistics with BNCweb – a Practical Guide*. Frankfurt a.M. [etc.]: Peter Lang. (= English Corpus Linguistics 6)
    - [CQP Query Language Tutorial](http://cwb.sourceforge.net/files/CQP_Tutorial.pdf) ([online version](http://cwb.sourceforge.net/files/CWB_Encoding_Tutorial/))
      http://cwb.sourceforge.net/files/CQP_Tutorial.pdf   http://cwb.sourceforge.net/files/CWB_Encoding_Tutorial/

# CQP-Abfragen: einzelne Tokens

- reg. Ausdruck in Anführungszeichen „matcht" Oberflächenform von Tokens
  - `"(over|under)\w+"` / `'(over|under)\w+'`
  - um Anführungszeichen zu finden: `""""` / `'"'`
  - immer regulärer Ausdruck ➞ Escapes bei Metazeichen nicht vergessen

- Suchoptionen:
  - `"deja"%c`     …     *case-insensitive* (Groß- und Kleinschreibung ignorieren)
  - `"deja"%d`     …     *diacritic-insensitive* (Diakritika ignorieren)
  - `"deja"%cd`     …     beides
  - `"?"%l`     …     exakte Zeichenkette (keine Metazeichen)

- Beispiele hier in BNCweb mit dem CLAWS C5-Tagset

# CQP-Abfragen: einzelne Tokens

- Abfage von Token-Annotation mit einem Attribut-Wert-Paar
  aus p-Attribut und regulärem Ausdruck für den Wert des Attributs:
  - `[lemma = "(over|under)\w+_ADJ"]`     (im BNC)
  - `[pos = "AJS"]`                          … Superlative (BNC)
  - `"deja"%cd` ist Kurzform für `[word = "deja"%cd]`

- Bedingungen mit Booleschen Operatoren kombinieren:
  - Operatoren: `&` (und), `|` (oder), `!` (nicht), `!=` (ungleich)

    **token description**

  - `[(word="can"%c) & (pos!="VM.*")]`
  - äquivalent zu `[(word="can"%c) & !(pos="VM.*")]`

# CQP-Abfragen: Token-Folgen

- CQP-Abfragen sind reguläre Ausdrücke über *token descriptions* (`[…]`)
  - `"in" [pos="AJ.*"]? [hw="year"]` … optional
  - `"in" [pos="AJ.*"]+ [hw="year"]` … mindestens eins
  - `"in" [pos="AJ.*"]{2} [hw="year"]` … genau zwei
  - `([pos="AJS"] | "most"%c [pos="AJ0"])` … entweder oder

- Abstände
  - `[]` … matchall (beliebiges Token)
  - `"dog" []{0,4} "cat"` … maximal vier Token dazwischen
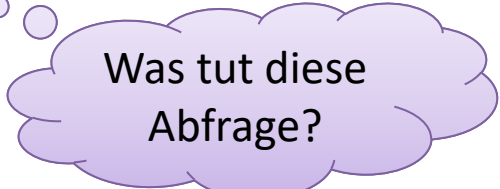  - `"dog" []{0,4} "cat" within s` … im selben Satz (s-Attribut)

# CQP-Abfragen: s-Attribute

- XML-Tags „matchen" Start- und End-Tag von s-Attribut-Bereichen (*regions*)
  - `<head> "UK"`          … als erstes Wort einer Überschrift
  - `"UK" </head>`          … als letztes Wort einer Überschrift
  - `<head>`                  … findet nichts (weil Länge von 0 Token)
  - `<mw> []* </mw>`    … passende Start-/End-Tags matchen ganzen Bereich

- innerhalb eines Bereichs suchen:
  - `"Twain" within quote;`
  - `[pos="NN.*"] :: match.mw_pos = "PRP";`
    … Abfrage von s-Attribut-Annotationen im „global constraint"
    … so können auch Bedingungen an Metadaten in der Query gestellt werden
  - vordefinierte Anker: `match`, `matchend`, `target` (`@`)

# CQP-Abfragen: Token-Folgen

- Quantoren und Alternativen können verschachtelt werden, um komplexe lexikalisch-grammatische Muster zu finden:

```
([pos="AJS"] | "most"%c [pos="AJ0"])
(
    "(and|\,)"%c
    ([pos="AJS"] | "most"%c [pos="AJ0"])
)+
[pos = "NN.*"]
```

> Was tut diese Abfrage?

- Voreinstellung: möglichst wenig Wiederholungen (*non-greedy*)
  - `"ho"%c ("," "ho"%c)+`          … findet immer *ho, ho*
  - `(?longest) "ho"%c ("," "ho"%c)+`
    … Strategie auf *greedy matching* umstellen (in CQPweb auch über UI möglich)

# Überlegungen vor der Recherche (1)

Vor einer ernsthaften Korpusrecherche gilt es, einige wichtige Überlegungen anzustellen:

- Was ist die sprachliche Grundgesamtheit, die einen interessiert? Standarddeutsch? (In Deutschland? In der Schweiz? Geschrieben und/oder gesprochen?) Eine bestimmte Varietät? Eine bestimmte Textsorte? Texte mit einem bestimmten Themenbezug? Texte aus den letzten paar Jahren oder aus vergangenen Jahrhunderten?
- Gibt es ein Korpus oder Teilkorpus, das dafür geeignet ist?
- Ist das gewählte Korpus groß genug? Ist daher für das zu untersuchende Phänomen mit genügend Treffern zu rechnen, um eine differenzierte Auswertung zu ermöglichen?

# Überlegungen vor der Recherche (2)

- Wie ist das Korpus genau aufgebaut? Welche Texte welcher Autoren, Textsorten, Entstehungszeiten usw. enthält es?

- Welche Zweifel an der Repräsentativität des Korpus gibt es? Wie gravierend sind diese? Worauf ist bei der Auswertung zu achten?

Wenn sich kein passendes Korpus finden lässt oder erste Testabfragen keine oder nicht genügend Treffer liefern, muss man u.U. darüber nachdenken, ein eigenes Korpus zusammenzustellen (oder die Fragestellung ganz anders anzugehen – z.B. mit einer Befragung).

# Übersicht: Web-Interfaces für Korpusabfragen

# CQPweb

- [https://corpora.linguistik.uni-erlangen.de/cqpweb/](https://corpora.linguistik.uni-erlangen.de/cqpweb/)
    - Login:        studentX (1 … 15)
    - Passwort: erlangen

CEQL*



- Hintergrund:
    - Hardie (2012); Evert & Hardie (2011)
    - [http://cwb.sourceforge.net/](http://cwb.sourceforge.net/)

- Dokumentation: Tutorial-Videos auf YouTube
  [https://www.youtube.com/user/CorpusWorkbench](https://www.youtube.com/user/CorpusWorkbench)

*http://cwb.sourceforge.net/ceql.php

# Sketch Engine

- https://app.sketchengine.eu/
  - sehr große Auswahl von Korpora in vielen verschiedenen Sprachen verfügbar
  - kommerzieller Anbieter: €90 / Jahr und Person (für Forschungszwecke)
  - Erstellen/Hochladen eigener Korpora bis zu einer Größe von 1 Mio. Wörter möglich (Upgrades möglich: z.B. auf 50 Mio. Wörter für €300 / Jahr)
  - FAU hat leider keine Uni-Lizenz (€3000 / Jahr)

- Bedienungsanleitung:
  https://www.sketchengine.eu/user-guide/

# Andere Web-Oberflächen @ FAU

- BNCweb (British National Corpus)
  https://corpora.linguistik.uni-erlangen.de/bncweb/
  - Login:        studentX
  - Passwort: erlangen
  - v.a. auch für Beispiele und Übungsaufgaben aus *Corpus Linguistics with BNCweb – a Practical Guide* (Hoffmann et al. 2008)

  CEQL*

- Debatten aus dem Europäischen Parlament
  http://corpora.linguistik.uni-erlangen.de/demos/CQP/Europarl/

  CEQL*

- deutsche Zeitungen aus 1990er Jahren (HGC: Huge German Corpus)
  http://corpora.linguistik.uni-erlangen.de/demos/auth/HGC/
  - Login:        demo
  - Passwort: demo
  - morphologisch annotiert

  CEQL*

# Weitere Web-Oberflächen (CWB-basiert)

- **OPUS**: Sammlung von Parallelkorpora
  http://opus.nlpl.eu/

- Leeds **IntelliText** (diverse Sprachen, Web-Korpora)
  http://corpus.leeds.ac.uk/itweb/htdocs/Query.html

- BFSU **CQPweb** (chinesische & englische Korpora @ **BFSU**)
  http://111.200.194.212/cqp/                          http://www.bfsu-corpus.org/channels/corpus

- **Linguateca AC/DC** (Portugiesisch)
  http://www.linguateca.pt/ACDC/

- **Ungarisches Nationalkorpus**
  http://corpus.nytud.hu/mnsz/index_eng.html

- Corpus del **Español Actual** (Spanisch)
  http://spanishfn.org/tools/cea/english

- **Varitext** (Französisch)
  http://syrah.uni-koeln.de/varitext

- **Spraakbanken** (Schwedisch)
  https://spraakbanken.gu.se/korp/

- **KorpusDK** (Dänisch)
  http://ordnet.dk/korpusdk/

# Weitere Web-Oberflächen (CWB-basiert)

- **TSCorpus** (Türkisch)
  http://tscorpus.com/

- **CORIS/CODIS** (Italienisch)
  http://corpora.ficlit.unibo.it/

- SSLMIT **La Repubblica** (italienische Zeitungen)
  http://dev.sslmit.unibo.it/corpora/corpus.php?path=&name=Repubblica

- **BwanaNet** (Katalanisch, Spanisch, Englisch)
  http://bwananet.iula.upf.edu/

- Georgetown University **CQPweb** (diverse Korpora)
  https://corpling.uis.georgetown.edu/cqp/

- **Perugia Corpus** (Italienisch)
  https://www.unistrapg.it/cqpweb/

- **CorpusEye** (diverse Sprachen)
  http://corp.hum.sdu.dk/

- **TEI:TOK** (mehrere Sprachen)
  http://www.teitok.org/index.php?action=projects

# Andere beliebte Web-Oberflächen

- [BYU Corpora](http://www.english-corpora.org/) (von Mark Davies)
  http://www.english-corpora.org/
  - COCA, COHA, Seifenopern, GloWbE, TIME, Spanisch, Portugiesisch, …

- Google [Web 1T 5-Grams](http://corpora.linguistik.uni-erlangen.de/cgi-bin/demos/Web1T5/Web1T5_freq.perl) (N-gramm-Datenbank)
  http://corpora.linguistik.uni-erlangen.de/cgi-bin/demos/Web1T5/Web1T5_freq.perl
  - N-Gramme durchsuchen, vorberechnete (Quasi-)Kollokationen
  - [NetSpeak](http://www.netspeak.org/): hübschere Web-Oberfläche
    http://www.netspeak.org/

- Google Books [Ngram Viewer](https://books.google.com/ngrams/) ([Info](https://books.google.com/ngrams/info))
  https://books.google.com/ngrams/  https://books.google.com/ngrams/info
  - Visualisierung von Frequenzänderungen im Laufe der Zeit (Wörter, Phrasen)

- Linguee: [Englisch](http://www.linguee.com/), [Deutsch](http://www.linguee.de/), [Französisch](http://www.linguee.fr/)
  http://www.linguee.com/  http://www.linguee.de/  http://www.linguee.fr/
  - Web-crawled parallel corpora for many language pairs
  - nützlich, um mögliche Übersetzungen zu finden (*caveat emptor* …)

- [Treebank.info](http://treebank.info/) (automatisch syntaktisch annotierte Korpora)
  http://treebank.info/

# Andere beliebte Web-Oberflächen

- **DWDS-Korpora** *
  https://www.dwds.de/r | http://kaskade.dwds.de/dstar/
  - Korpusabfragen, (diachrone) Kollokationsanalyse, Wortverlaufskurven

- **COSMAS II** *
  https://cosmas2.ids-mannheim.de/cosmas2-web/
  - Korpora des Instituts für deutsche Sprache (IDS), inkl. Kollokationsanalyse

- **DGD**: Datenbank für gesprochenes Deutsch des IDS *
  https://dgd.ids-mannheim.de

- **KorAP**: neue Web-Oberfläche für IDS-Korpora (CQL-Abfragen möglich)
  https://korap.ids-mannheim.de/

- **ANNIS-Interface** der HU Berlin: diverse, v.a. kleinere Korpora
  https://korpling.german.hu-berlin.de/annis3/

\* Kurzanleitung für DWDS, COSMAS II und DGD:
http://sprachwissenschaft.fau.de/personen/daten/blombach/korpora.pdf