

# Korpuslinguistik I

Morphologieanalyse und Lexikonbau (8. Vorlesung)

Dozent: Gerold Schneider

## Übersicht

- [Was ist Korpuslinguistik?](#)
- [Anwendungen der Korpuslinguistik](#)
- [Grundfragen der Korpuslinguistik](#)
- [Korpusaufbau](#)
- [Ein Konkordanzprogramm im Überblick: Conc 1.8 für Macintosh](#)
- [Beispiel einer Visualisierung in MS Excel](#)

## Literatur:

- [\[McEnery 96\]](#): Immer noch das klassische Einführungsbuch
- [\[Biber 98\]](#): Ausführliches Kapitel über Lexikographie (Part I Chapter 2)
- [\[Kennedy 98\]](#): Gute Übersicht über englische Korpora (Chapter 2)

# Was ist Korpuslinguistik?

Korpuslinguistik (im folgenden KL)

- ist Linguistik basierend auf einer grossen Textsammlung=Korpus
- erhielt ihren Namen von einem Instrument der Analyse (wie etwa *Computerlinguistik*)
- hat die gesamte Linguistik revolutioniert
- verwendet statistische Methoden über quantitativ-linguistischen Analysen

## Korpuslinguistik (KL) und Computerlinguistik (CL)

KL wird von vielen Computerlinguisten nur als Randgebiet der CL aufgefasst, aber:

- Viele CL Anwendungen basieren auf KL
- KL liefert statistische Informationen für CL (z.B. für Disambiguierung)
- KL ist ein Bindeglied zwischen klassischer Linguistik und CL
- Symbolische versus subsymbolische Ansätze (Z.B. KI versus Statistik)
- KL ist prinzipiell auch ohne Computer möglich, obwohl dann extrem aufwendig

## **Präskriptive versus deskriptive Linguistik**

Wörterbücher, Schul- und Fremdsprachengrammatiken werden häufig zu Rate gezogen, um den richtigen Gebrauch der Sprache sicherzustellen. Die Linguistik vergangener Jahrhunderte sah einen ihrer Hauptzwecke darin, den richtigen Gebrauch der Sprache zu kontrollieren: Präskriptive Linguistik.

Die deskriptive Linguistik hingegen beschreibt Sprache so, wie sie wirklich auftritt.

## **Performanz versus Kompetenz**

Siehe [\[McEneaney 96\]](#) Chapter 1.

Unter Sprachkompetenz versteht man die Fähigkeit eines Sprechers, wohlgeformte Sätze aufgrund von Sprachregeln zu bilden. Chomsky spricht auch von I(nternalisierter) Sprache und versteht darunter ein Regelsystem (z. B. Phrasenstrukturregeln und Transformationen), das alle grammatischen Sätze generieren kann und alle ungrammatischen Sätze zurückweist.

Unter Performanz versteht man die Summe der Sprachäußerungen, die ein Sprecher einer Sprache von sich gibt. Chomsky spricht auch von E(xternalisierter) Sprache, eine extrem lange Auflistung aller Sätze, die in irgendeinem Kontext je geäußert wurden.

Kompetenz wird auch als Langue, Performanz auch als Parole bezeichnet.

Chomsky sah die Aufgabe der Linguistik vor allem darin, ein psycholinguistisches Modell der Sprachkompetenz zu bauen und verurteilte Performanz als einen untauglichen Spiegel der Kompetenz.

Zwischen Mitte der 50er und den 80er Jahren kam die Forschung in der KL fast zum Aber das Chomskyanische Universalgrammatikmodell liegt immer noch in endloser Ferne liegt und man kann Sprachdaten sinnvoll einsetzen kann

- zur Überprüfung von Hypothesen
- zur Anreicherung von Regeln mit statistischen Wahrscheinlichkeiten, was Ambiguität verringert
- zur detaillierten Sprachbeschreibung, viel zuverlässiger als die Intuition eines Grammatikschreibers
- als Grundlage für eine empirisch-wissenschaftliche Beschreibung der Sprache

→ Sprachwissenschaft an der Schnittstelle zwischen einer empirischen und einer rationalen Wissenschaft.

KL ist das systematische, streng empirische Vorgehen mit Beobachtung und Regelableitung

Der introspektive ("armchair") Linguist geht dagegen in der Art des Rationalismus vor

# Anwendungen der Korpuslinguistik

KL, da eher linguistisches Instrument als Teilgebiet der Linguistik, findet Anwendungen in fast allen linguistischen Teilgebieten

## Morphologie

- Produktivität: Auf der Basis eines morphologisch annotierten Korpus oder mittels eines zuverlässigen (!) Morphologieanalyseprogramms: Welche Derivationsmorpheme sind produktiv? Welche Partikel- und Derivationsmorpheme treten bei welchen Kernmorphemen auf, welche bei fast allen (produktiv), welche nur sehr eingeschränkt?
- Flexion: Welche Flexionsformen treten häufig, welche selten auf (sg/pl,nom/acc, 1./2./3. Person) und gibt es Registerunterschiede? Stimmen unsere Annahmen über Markiertheit und Ikonizität? Werden unzählbare Substantive wirklich nie im Plural verwendet?
- Tagging: Welche Wortart von allen möglichen einer gegebenen Wortform ist häufiger? Und in welchem Kontext? Kann die richtige Wortart aufgrund dieser Angaben immer zugewiesen werden? (siehe [Statistisches Tagging](#) und [Regelbasiertes Tagging](#))
- Distribution: Welche Worte und Wortarten sind wie verteilt? Treten gewisse Wortarten in gewissen Sprachen oder Registern häufiger auf?

## Lexikologie und Lexigraphie

[\[Biber 98\]](#) Chapter 2 ist ganz diesem Thema gewidmet.

Lexikographen sind schon seit Jahrhunderten bemüht, auf KL abzustützen. Samuel Johnson suchte schon 1755 mühevoll nach Beispielen aus der Literatur für die Einträge in seinem Lexikon. Heute gibt es kaum mehr ein Lexikon, das nicht KL-getestet ist.

KL wird heute verwendet für:

- Semi-automatische Lexikographie
- Semi-automatische Term-Extraktion: Z.B.: Tritt ein best. Wort gehäuft oder nur in einem gewissen Register auf, oder fehlt es in einem allgemeinen Lexikon, so ist es ein Term-Kandidat.

KL-basierte Methoden erlauben zusätzlich lexikologische Angaben zu:

- Sprachgebrauch: Welche deutschen Worte sind die häufigsten (Strukturworte)? Welche Bedeutung eines polysemischen Wortes kommt wie häufig vor?
- Kollokationen: In welchen Kontexten (lokal oder hierarchisch) treten welche Worte auf? Sind diese für Synonyme gleich? Welche grammatischen (Syntax) und semantischen Begleiter (semantische Assoziationen) treten auf?

## Syntax

Gegeben sei folgende Trivialgrammatik und ein Parsingaufruf:

```
S -> NP VP.  
S -> NP.  
NP -> N N.  
NP -> N.  
VP -> 'walks'.  
NP -> 'walks'.  
N -> 'John'.
```

Parse>> John walks.

Wie könnte man mit Hilfe eines annotierten Korpus der richtigen Lesart dieses Satzes Vorzug geben?

Häufig werden statistische Gewichtungen zur Disambiguierung von PP-Anbindung oder Adverbskopos verwendet:

John sees the man in the park with the telescope.

Welche Verben und welche Nomen verwenden welche Präpositionen mit Vorliebe?

Heureusement, Jean part pour les vacances.  
Jean heureusement survit l'accident.

Welche Verben und Adverben (oder Kombinationen) haben eher welchen Skopus?

## **Semantik**

In der Semantik gibt es z.B. Untersuchungen darüber, ob zwei synonymische Ausdrücke tendenziell von den gleichen Worten begleitet werden (Kollokationen). Es gibt auch Versuche, semantische Netzwerke zu konstruieren, bei denen (für referenzierende Worte) die durchschnittliche lokale Nähe im Text als Parameter der semantischen Nähe angenommen wird.

## **Phonetik**

Ein grosser Teil der phonetischen KL untersucht Sprachvariation nach Alter, Dialekt, Ausbildung, Geschlecht etc. Korpora mit gesprochenen Daten sind aber auch interessant da sie noch näher an 'real-world' Sprachgebrauch sind als Textcorpora, also z.B. Grammatik so, wie sie 'in unseren Köpfen' lebt anstatt in Schulbüchern. Dies ist für die Psycholinguistik von Interesse. Viel linguistisches Wissen über Prosodie wurde auch durch Korpusanalysen gewonnen. Grosse Corpora gesprochener Sprache werden in der CL verwendet, um Spracherkennungssysteme möglichst sprecherunabhängig trainieren zu können.

## **Pragmatik, Diskursanalyse**

Meist auch anhand gesprochener Korpora. Wer spricht wann, wie sind die Übergänge zwischen Sprechern, wer unterbricht eine Äusserung, wie signalisieren Hörer ihr Interesse an gesprochenen?

## **Soziolinguistik, Dialektologie**

Soziolinguistik ist per definition deskriptive Linguistik. Korpusbasierte Ansätze untersuchen Sprachvariation nach Alter, Dialekt, Ausbildung, Geschlecht und so weiter.

## **Historische Linguistik**

Den meisten qualitativen Veränderungen in der Sprache geht ein langer schleichender soziolinguistischer Prozess voran, in dem eine seltenere Form die üblichere immer stärker vertritt und diese mit der Zeit überflügelt. Nur quantitative Untersuchungen belegen die Präferenzen bei der Auswahl, die letztlich zu einer qualitativen Änderung führen.

## **Stilistik und Stilometrie**

Verschiedene Autoren bedienen sich verschiedener Stilmittel, was sich auch in der quantitativen Sprachbeschreibung niederschlägt. Welcher Autor braucht welche Worte häufiger, macht typischerweise welche Fehler, hat ungefähr welchen Substantiv/Verb Koeffizient etc. Die Stilometrie kann häufig Autoren unbekannter Werke identifizieren oder stellt sich die Frage, ob Shakespeare nur eine Person oder das Pseudonym einer Autorengruppe war.

# Grundfragen der Korpuslinguistik

## **Zusammenstellung des Korpus: Ausgewogenheit, Repräsentativität**

Die Auswahl der Texte, die in einen Korpus aufgenommen werden sollen, ist entscheidend. Will man einen genrespezifischen oder einen allgemeinen Korpus zusammenstellen? Welche Texte repräsentieren welches Genre? Schwierig: Wie sollen verschiedene Genres in einem allgemeinen Korpus vertreten sein? Letztlich gibt es prinzipiell nie einen völlig ausgewogenen Korpus.

## **Statistische Methoden: Aussagekraft, Relevanz**

Alleine die Erkenntnis, dass ein gewisses Phänomen in einem Korpus häufiger als in einem anderen vorkommt, macht noch keine statistisch aussagekräftige Aussage. Eine dialektologische Untersuchung, die nur auf einen Probanden jeder Dialektgruppe abstützt, ist statistisch vollkommen wertlos, da die Unterschiede genausogut individueller Präferenzen anstatt dialektaler entspringen können.

Falls der zu untersuchende Häufigkeitsunterschied zwischen zwei (Teil-)Korpora gering ist oder falls nur eine kleine Anzahl Belege gefunden werden konnten, so ist das Risiko gross, dass es sich beim Unterschied nur um zufällige Fluktuationen handelt. Um den Einfluss zufälliger

Schwankungen einschätzen zu können und somit echte von zufälligen Unterschieden trennen zu können, verwendet man eine Reihe von statistischen Tests. Siehe [nächste Vorlesung](#).

## **Korpusannotation**

Ein roher Text erlaubt nur klare Aussagen auf der Wortformenebene. Für gültige Aussagen auf der phonologischen, morphologischen, syntaktischen oder semantischen Ebene darf man sich nicht exklusiv auf fehlerbehaftete automatische Analysemethoden abstützen. Auf der morphologischen Ebene beispielsweise macht auch eine Taggingfehlerrate von nur 2% u.U. eine statistische Aussage völlig zunichte, falls sie teilweise auf die fehlerhaften Tags in einem so getaggen Korpus abstützt.

Deshalb ist es meist notwendig, sich auf zuverlässige, handkorrigierte oder -erstellte Analysen abzustützen, mit denen der Korpus annotiert ist. Eine Annotation auf möglichst vielen Ebenen ist zwar wünschenswert, verteuert und verlängert aber die Zusammenstellung eines Korpus ins Endlose. So beschränkt man sich meist auf wenige Ebenen oder der Korpus bleibt sehr klein. Für viele syntaktisch interessante Aussagen sind die meisten (aufwendig!) syntaktisch annotierten Korpora auch heute noch zu klein.

Als Beispiel einer morphosyntaktischen Annotation siehe [STTS aus der ersten Vorlesung](#).

# Verwaltung der Korpusdaten: Speicherung, Zugriff, Visualisierung

## Speicherung

Die Verwaltung von grossen Datenmengen stellt Probleme in der Art der Speicherung, die schon auf einen einfachen und schnellen Zugriff vorbereiten soll.

- **Datenbank (DB):** Effizient, auf schnellen Zugriff ausgelegt. Die Textdaten müssen aber konvertiert werden in ein DB-Format, was aufwendig und für Linguisten nicht unbedingt einfach ist. Die konvertierten Daten sind für Linguisten nicht mehr manuell lesbar.
- **Rohtext:** Sehr flexibel, der Quelltext kann direkt verwendet werden, keine DB-Spezialisten notwendig. Die Speicherung kann aber platzintensiv und der Zugriff langsam sein. Es ist schwierig, so einen Korpora zu handhaben, der nicht vollständig im RAM Platz hat.
- **SGML, XML:** Strukturierter Text, schneller als Rohtext, langsamer als DB, aber manuell lesbar. Ein Kompromiss zwischen DB und Rohtext.
- **proprietäre Formate:** Auch effizient, aber schwierig portierbar und überhaupt nicht mehr manuell lesbar.

## Zugriff

- **SQL** (oder andere Datenbankabfragesprache): Der übliche Zugriff auf eine Datenbank.  
Sehr effizient, ziemlich flexibel.

Ein SQL-Abfragebefehl hat das folgende Format:

```
SELECT Zu findende Werte
FROM Tabelle
WHERE Bedingung
```

Stellen wir uns z.B. folgenden einfachen Korpusaufbau vor:

...

SATZ 2112	1	2	3
wort	Ein	Satz	!
tag	DET	NN	PKT

SATZ 2113	1	2	3	4	5	6	7
wort	Der	vorliegende	Text	ist	ein	Beispiel	.
tag	DET	ADJA	NN	VVFIN	DET	NN	PKT

SATZ 2114	1	2	3	4	5
wort	Hier	geht	es	weiter	.
tag	ADV	VVFIN	PRO	PTKVZ	PKT

...

Dann würde beispielsweise eine Anfrage nach den Artikeln im 2113. Satz des Korpus folgendes Aussehen haben, und liefert die darunterstehenden Ergebnisse:

```
SELECT wort
FROM SATZ 2113
WHERE tag = DET
```

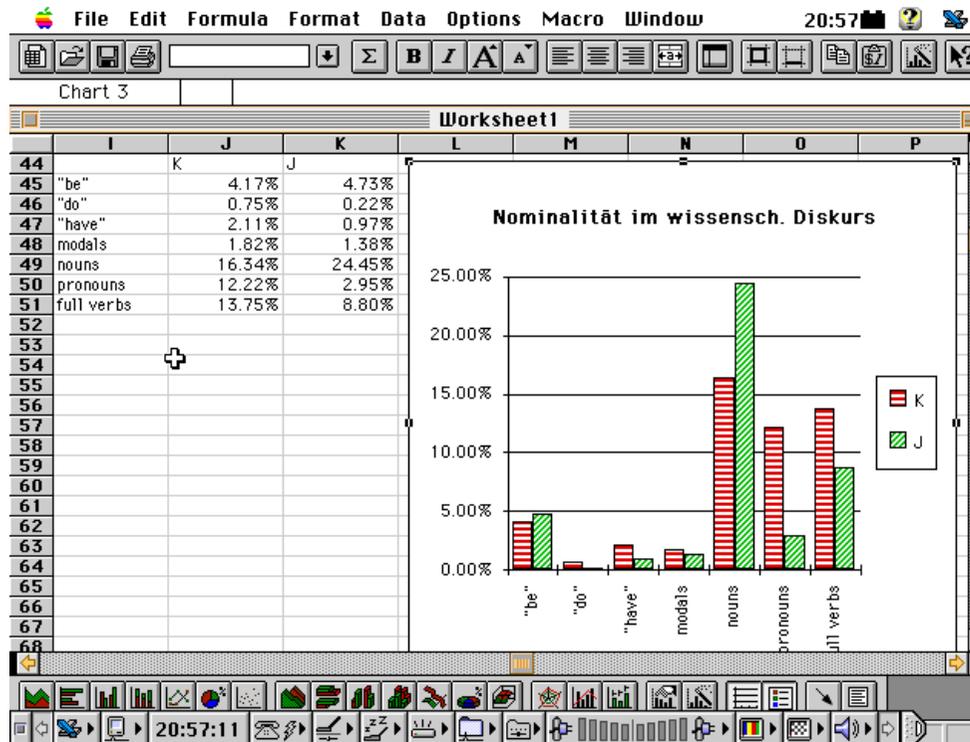
```
> wort = Der
> wort = ein
> ALL RESULTS
```

- **Konkordanzprogramme:** Einfach, oft visuelle Ausgabe der Ergebnisse. Die Möglichkeiten sind aber beschränkt und nicht erweiterbar. Siehe [ein Beispiel unten](#).
- **Perl** (oder andere listenoptimierte Programmiersprache): Arbeitet über Rohtext, Listenvariablen oder gar DB-Formate. Eher langsam, aber extrem flexibel.

```
#!/local/bin/perl
# Dieses Programm sucht STTS Artikel, gibt sie aus und z"ahlt die Types
# (c) GS
while (<>) {
    ($word, $tag) = split(/\_\/);
    if ($tag == "DET") {
        print "$word \n";
        ++$typecounter{$word};
    }
}
# lese Korpus wortweise ein
# trenne Wort und Tag bei _
# ist es ein Artikel (DET)? Falls ja:
# gebe Wort aus
# erhöhe Zähler für dieses Wort um 1
# Ende der while-Schleife: Nimm nächstes Wort
# und durchlaufe mit ihm die Schleife, etc.
```

## Visualisierung

Die extrahierten Daten, in der Regel lange Listen und Tabellen, geben oft noch keinen Überblick über die untersuchten Zusammenhänge. Hier empfiehlt sich der Einsatz eines professionellen Statistikprogrammes oder (für einfachere Aufgaben) des altbekannten MS Excel.



# Korpusaufbau

Ein historischer wie thematischer Abriss über den Aufbau der ersten grösseren computerlesbaren Korpora.

## Aufteilung der Texte

Schon 1964 begann in den USA an der Brown Universität (daher der Name Brown Corpus) die Zusammenstellung des ersten grösseren computerlesbaren Corpus. Die Zeiten waren schwierig; "it was compiled in the face of massive indifference if not outright hostility from those who espoused the conventional wisdom of the new and increasingly dominant paradigm in US linguistics led by Noam Chomsky" ([\[Kennedy 98\]](#) p. 23). Die damals immense angepeilte Korpusgrösse von 1 Mio Worte wurde innerhalb von 3 Jahren Arbeit erreicht. Im Bestreben um Ausgewogenheit wurden folgende Genres aufgenommen:

### Aufteilung des Brown und des LOB Korpus

Category A (Press: reportage)  
Category B (Press: editorial)  
Category C (Press: reviews)  
Category D (Religion)  
Category E (Skills, trades and hobbies)  
Category F (Popular lore)  
Category G (Belles lettres, biography, essays)

Category H (Miscellaneous, mainly Government documents)  
Category J (Learned and scientific writings)  
Category K (General Fiction)  
Category L (Mystery and detective fiction)  
Category M (Science fiction)  
Category N (Adventure and western fiction)  
Category P (Romance and love story)  
Category R (Humour)

Diesseits des Atlantik standen erst in den 70er Jahren genügend Ressourcen zur Verfügung um ein ähnliches Projekt in Angriff zu nehmen. Von 1970-78 wurde als gemeinsames Projekt der beteiligten Universitäten der Lancaster-Oslo-Bergen (LOB) Korpus zusammengestellt, nach der gleichen Aufteilung wie der Brown Korpus, ebenfalls mit etwa 1 Mio Worten, womit beispielsweise erstmals ein quantitativer Vergleich zwischen US und GB Englisch möglich wurde.

1988 folgte ein fast gleich aufgeteilter Korpus für Indisches Englisch, darauf ebenso einer für Australisches, und 1993 einer für Neuseeländisches Englisch. Auch heute noch finden diese Korpora weite Verwendung, nicht zuletzt weil sie klein genug sind um auf heute üblichen PCs bearbeitet zu werden.

## **Annotierung der Korpora**

Schon in den 70er Jahren wurden erste Versuche zum semi-automatischen part-of-speech Tagging vorgenommen. Die Fehlerraten waren aber noch sehr hoch. Erst das für den LOB

Korpus entwickelte CLAWS System erreichte Fehlerraten von unter 5%. Die manuelle Korrektur bleibt aufwendig, im Falle das LOB Korpus erfolgte sie zwischen 1978 und 1983.

Teile dieser Korpora liegen auch schon syntaktisch annotiert in hierarchischer Phrasenstruktur vor, aber noch keiner von ihnen vollständig. Parsing und manuelle Korrektur sind noch wesentlich aufwendiger als fürs part-of-speech Tagging.

## Heutige Korpora

Heute gibt es viele Spezialkorpora in Englisch, für gesprochenes Englisch, diachrone Sprachentwicklung, Spracherwerb etc. Auch fürs Deutsche sind schon einige erhältlich (z.B. der STTS-annotierte Frankfurter Rundschau Korpus mit etwa 250'000 Worten oder der kleine von uns STTS-annotierte [Zürcher Universitätskorpus](#) mit 50'000 Worten).

Im Englischen wird die Entwicklung von extrem grossen Corpora vorangetrieben. Der aus dem Cobuild-Projekt entstandene, allerdings ungeaggte Bank of English Korpus umfasst etwa 300 Mio Worte. Der British National Corpus (BNC) ist ähnlich wie Brown oder LOB getaggt und umfasst 100 Mio Worte, aber es ist fraglich, wieviele Taggingfehler die gezwungenermassen oberflächliche Korrektur des automatischen Tagging noch überlebt haben. Am Englischen Seminar der Uni Zürich ist [Prof. Totties Gruppe](#) an der BNC-Forschung vertreten.

[\[Kennedy 98\]](#) Chapter 2.3 bietet einen Überblick über einige weitere Korpora.

# Ein Konkordanzprogramm im Überblick: Conc 1.8 für Macintosh

Ein Konkordanzprogramm erlaubt

- das Auffinden aller Vorkommnisse eines Suchwortes im Kontext (KWIK)
- das Erstellen von sortierten Wort- und Annotationslisten
- Kollokationsanalysen
- einfache Statistik und tabellarische Ausgabe von Daten für die Weiterverarbeitung

Conc (derzeitige Version ist 1.80b3) kann man [frei herunterladen](http://www.sil.org) vom amerikanischen Summer Institute of Linguistics ([www.sil.org](http://www.sil.org)). Es ist relativ einfach, sowohl in der Bedienung als auch in den Möglichkeiten. Es verwendet z.B. Korpora in rohem Textformat, was maximal einfach ist, aber lange Verarbeitungszeiten zum Aufbau der Konkordanz mit sich zieht. Ein Nachteil der Einfachheit ist auch, dass der ganze Korpus im RAM Platz finden muss.

File Edit Font Options Layout Build Windows 15:00

LOBTH\_K.THT

K26 46 lunch\_NN with\_IN me\_PP10 yesterday\_NR ... \*\*'\_\*\*'^ notes\_NNS  
 K26 47 would\_MD leave\_VB saying\_VBG ^ '\*\_\*' sorry\_JJ , my\_PP\$ tutor\_NN  
 K26 47 came\_VBD ... ^ ( ( my\_PP\$ tutor\_NN would\_MD often\_RB pop\_VB  
 K26 48 in\_RP , and\_CC we\_PP1 AS would\_MD retire\_VB to\_IN a\_AT nearby\_JJB  
 K26 48 teashop\_NN , eat\_VB buns\_NNS , and\_CC discuss\_VB my\_PP\$  
 K26 49 thesis\_NN , at\_IN the\_ATI same\_AP time\_NN feeding\_VBG crumbs\_NNS  
 K26 49 to\_IN the\_ATI mice\_NNS that\_WPR kept\_VBD  
 K26 50 appearing\_VBG out\_RP of\_IN the\_ATI wainscoting\_NN ... ) ) ^ but\_CC  
 K26 50 how\_WRB about\_IN today\_NR ?\_? ^ !\_PP1 A 'm\_BEM your\_PP\$

Concordance

10163 59 it\_PP3 was\_BEDZ less\_QL **easy\_JJ** to\_TO believe\_VB and\_CC which\_WDTR K29 59 Kungo\_NP had\_HVD never\_RB  
 7802 , K22 148 there\_RN to\_TO **eat\_VB** Wiener\_NP schnitzel\_NN at\_IN a\_AT long\_JJ K22 148 table\_NN with\_IN  
 7838 order\_TO" to\_TO" K22 166 **eat\_VB** , and\_CC stop\_VB eating\_VBG in\_TO order\_TO" K22 167 to\_TO" talk\_VB ,  
 8548 you\_PP2 do\_DO n't\_XNOT **eat\_VB** sweets\_NNS ?\_? K24 185 ^ \*\_\*\_ no\_UH ... K24 186 ^ \*\_\*\_ or\_CC a\_AT  
 9054 K26 48 teashop\_NN , **eat\_VB** buns\_NNS , and\_CC discuss\_VB my\_PP\$ K26 49 thesis\_NN , at\_IN the\_ATI  
 9901 his\_PP\$ hand\_NN ... ^ '\*\_\*' **eat\_VB** K28 131 your\_PP\$ dinner\_NN , darling\_NN ... '\*\_\*'\* K28 132 ^  
 10100 it\_PP3 did\_DOD not\_XNOT **eat\_VB** their\_PP\$ herds\_NNS K29 27 and\_CC they\_PP3AS could\_MD now\_RN drive\_VB  
 5470 ^ she\_PP3A K16 38 had\_HVD **eaten\_VBN** too\_QL much\_AP of\_IN the\_ATI K16 39 smoked\_JJ salmon\_NN at\_IN  
 7621 had\_HVD K29 53 your\_PP\$ **eat\_VB** K29 53 ^ \*\_\*\_ or\_CC in\_ATI NN ...

Index

easier_JJR	(3)	5034, 5038, 5417
easily_RB	(5)	2691, 6751, 7109, 8410, 10300
east_NR	(6)	410, 508, 515, 559, 4197, 7903
easy_JJ	(8)	717, 875, 3129, 3631, 7003, 7531, 9381, 10163
<b>eat_VB</b>	<b>(6)</b>	<b>7802, 7838, 8548, 9054, 9901, 10100</b>
eaten_VBN	(2)	5470, 7621
eating_NN	(1)	7715
eating_VBG	(4)	204, 7838, 7930, 9766

15:00:44

[Conc-Konkordanz mit LOB Korpus Kategorie K, Sortierfolge ASCII]

**LOBTH\_K.THT**

K01 6 pretty\_RB chipper\_JJ , thanks\_NNS , considering\_RI ... \*\_\*\*^  
 K01 6 he\_PP3A  
 K01 7 was\_BEDZ a\_AT tiny\_JJ man\_NN , of\_IN fanatical\_JJ neatness\_NN ,  
 K01 7 his\_PP\$ remaining\_JJ hair\_NN snowy\_JJ , and\_CC  
 K01 8 cropped\_VBN like\_IN a\_AT Prussian's\_NNP\$ ... ^ his\_PP\$ white\_JJ  
 K01 8 shirt\_NN cuffs\_NNS were\_BED actually\_RB  
 K01 9 starched\_VBN : he\_PP3A protruded\_VBD from\_IN them\_PP3OS his\_PP\$  
 K01 9 surprisingly\_RB thick\_JJ and\_CC hairy\_JJ  
 K01 10 wrists\_NNS and\_CC began\_VBD to\_TO wash\_VB ... ^ \*\_\*\*^ as\_IN a\_AT

**Concordance**

10256	careful_JJ in_IN obeying	_VBG K29 107 the_ATI custom_NN of_IN the_ATI tribe_NN ... ^ he_PP3A K29 108
10289	tired_JJ of_IN looking	_VBG at_IN K29 124 the_ATI wild_JJ animals_NNS ... ^ Bwana_NPT Dillon_NP
10300	white_JJ men_NNS talking	_VBG easily_RB and_CC K29 129 casually_RB of_IN the_ATI women_NNS with_IN K29
13	, and_CC K01 8 cropped	_VBN like_IN a_AT Prussian's_NNP\$ ... ^ his_PP\$ white_JJ K01 8 shirt_NN cuffs_NNS
15	actually_RB K01 9 starched	_VBN : he_PP3A protruded_VBD from_IN them_PP3OS his_PP\$ K01 9
24	a_AT voice_NN K01 13 made	_VBN resonant_JJ by_IN the_ATI very_AP K01 14 weakness_NN of_IN his_PP\$
50	than_RB" once_RB read	_VBN his_PP\$ K01 28 entry_NN in_IN who_WP 's_BEZ who_WP : son_NN of_IN
52	28 Fairbanks_NP , married	_VBN to_IN the_ATI K01 29 daughter_NN of_IN a_AT knight_NN , member_NN of_IN
57	K01 30 doubtless_RB	_VBN : he_PP3A protruded_VBD from_IN them_PP3OS his_PP\$ K01 9

**Index**

_TO	(1076)	3, 17, 19, 26, 45, 64, 116, 131, 168, 187, 205, 212, 217, 222, 225, 253, 324, 330, 3
_JH	(203)	41, 245, 245, 327, 386, 386, 429, 429, 445, 473, 479, 737, 795, 806, 809, 855, 882
_VB	(2484)	3, 17, 19, 26, 27, 32, 35, 35, 45, 61, 64, 65, 66, 72, 80, 94, 107, 116, 131, 168, 17
_VBD	(2820)	3, 5, 8, 15, 17, 21, 23, 26, 31, 33, 42, 54, 61, 88, 89, 95, 95, 96, 103, 111, 147, 15
_VBG	(987)	28, 40, 40, 112, 145, 160, 174, 186, 190, 200, 202, 204, 218, 250, 271, 294, 301, 3
_VBN	(1375)	13, 15, 24, 50, 52, 67, 69, 71, 76, 87, 91, 92, 99, 103, 106, 109, 119, 123, 126, 12
_VBZ	(151)	19, 220, 369, 377, 597, 741, 749, 796, 798, 910, 910, 911, 1014, 1078, 1266, 1418
_WDT	(201)	40, 103, 222, 233, 506, 532, 568, 657, 822, 936, 1021, 1073, 1084, 1102, 1209, 12

15:09:18

[Conc-Konkordanz mit LOB Korpus Kategorie K, sortiert nach Tags]

Einfache quantitative Untersuchungen kann man schon in Conc selber vornehmen. So z.B.:

## **Untersuchungsfrage: Ist wissenschaftlicher Diskurs nominaler als andere Genres?**

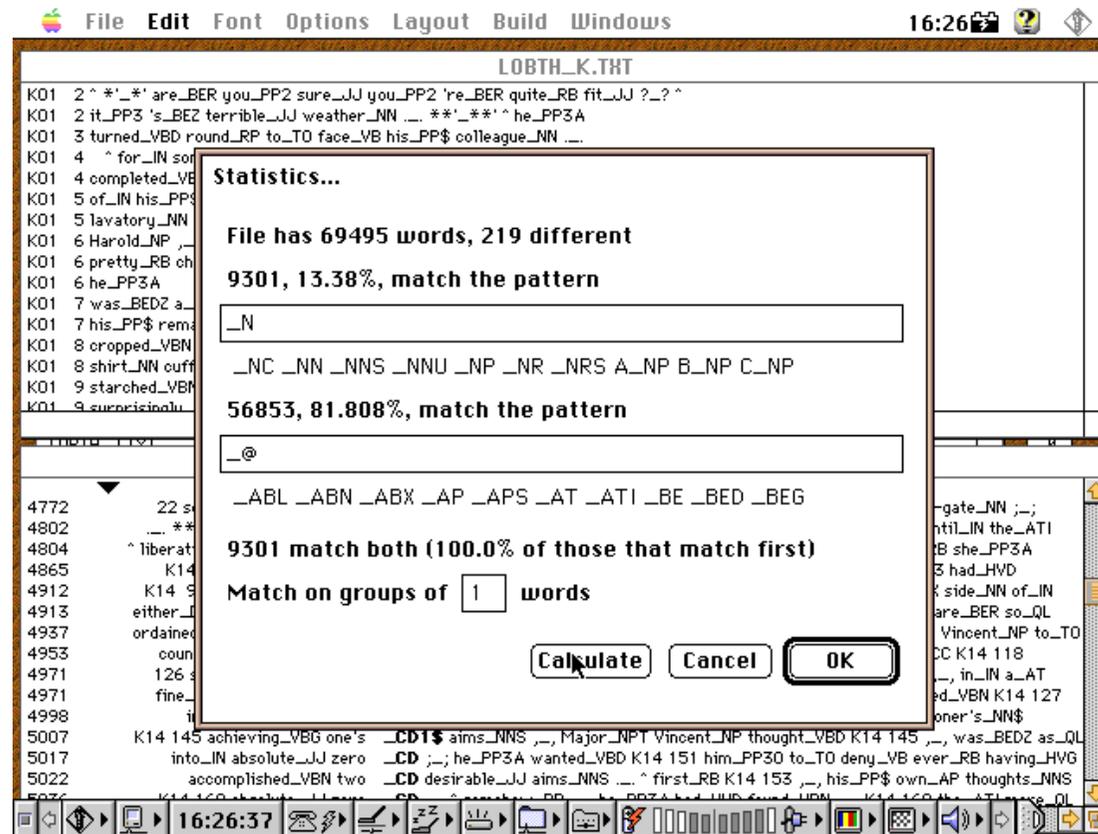
Wir beschränken uns auf einen Direktvergleich der LOB Kategorien K (General Fiction) und J (Learned and Scientific Writing) und versuchen, die folgenden einfachen Fragen zu beantworten:

- Wieviele Prozent aller Worte sind Nomen?
- Wieviele Prozent aller Worte sind Verben?
- Was ist das Verb/Nomen Verhältnis?
- Wie sind die Unterschiede zwischen LOB K und LOB J betreffend dieser Fragen?

→ Zur on-line Demo

Grob gesagt beginnen im LOB Tagset Verbtags mit V und Nomentags mit N. Die einfache in Conc eingebaute Statistikabfrage über reguläre Ausdrücke liefert uns dazu folgende Ergebnisse:

- Wieviele Prozent aller Worte sind Nomen?



## In der online-Demo werden wir

- Einen spezialisierten Index aufbauen

The screenshot shows the LOB-LexMorph software interface. The 'Include Words Options' dialog box is open, showing the following settings:

- Include all words
- Select words to include:
  - Include groups of up to 1 words that match one of these patterns:
  - Pattern:
  - Pattern:

The main window displays a concordance table with columns for line numbers, source text, and target text. Below the concordance is an 'Index' window showing a list of morphological tags and their frequencies:

Tag	Frequency	Line Numbers
_NNS	(2055)	6, 9, 14, 17, 19, 23, 31, 55, 57, 60, 61, 68, 71, 87, 96, 100, 105, 107, 110, 111, 112
_NNU	(9)	2361, 2361, 2364, 2364, 2367, 2367, 8481, 8771, 8771
_NP	(2)	9190, 9251
_NR	(79)	336, 353, 353, 363, 410, 460, 508, 515, 559, 595, 805, 874, 912, 1027, 1384, 1543
_NRS	(1)	1536
_VB	(2484)	3, 17, 19, 26, 27, 32, 35, 35, 45, 61, 64, 65, 66, 72, 80, 94, 107, 116, 131, 168, 17
_VBD	(2820)	3, 5, 8, 15, 17, 21, 23, 26, 31, 33, 42, 54, 61, 88, 89, 95, 95, 96, 103, 111, 147, 15
_VBG	(987)	28, 40, 40, 112, 145, 160, 174, 186, 190, 200, 202, 204, 218, 250, 271, 294, 301, 3

- Diesen als Rohtext exportieren

The screenshot shows the LOB-LexMorph software interface. The 'Include Words Options' dialog box is open, showing the following settings:

- Include all words
- Select words to include:
  - Include groups of up to 1 words that match one of these patterns:
  - Pattern:
  - Pattern:

The main window displays a concordance table with columns for line numbers, source text, and target text. Below the concordance is an 'Index' window showing a list of morphological tags and their frequencies:

Tag	Frequency
_BE	(294)
_BED	(274)
_BN	(1000)
_BZ	(46)
_CC	(104)
_CD	(172)
_CD1	(155)
_CD2	(328)
_CD3	(223)
_CD4	(167)
_CD5	(37)
_CD6	(279)
_CD7	(791)
_CD8	(29)
_CD9	(21)
_CD10	(81)
_CD11	(1033)
_CD12	(12)
_CD13	(7130)
_CD14	(2055)
_CD15	(9)
_CD16	(2)
_CD17	(79)
_CD18	(1)
_CD19	(249)
_CD20	(3180)
_CD21	(1761)
_CD22	(531)
_CD23	(124)
_CD24	(48)
_CD25	(819)
_CD26	(234)
_CD27	(2484)
_CD28	(2820)

A file export dialog box is open, showing the following settings:

- Save as:
- Buttons: Eject, Desktop, New, Cancel, Save

• In Excel tabellarisch verarbeiten

	E	F	G	H	I	J	K	L
14	343	_HYD	1.39%	0.22%				
15	64	_HYG	0.05%	0.04%				
16	23	_HYN	0.04%	0.01%				
17	469	_HVZ	0.14%	0.30%	"have"	2.11%	0.97%	
18	2135	_MD	1.82%	1.38%	modals	1.82%	1.38%	
19	56	_NC	0.02%	0.04%				
20	26933	_NN	12.54%	17.41%				
21	9897	_NNS	3.61%	6.40%				
22	775	_NNU	0.02%	0.50%				
23	9	_NP	0.00%	0.01%				
24	147	_NR	0.14%	0.10%				
25	5	_NRS	0.00%	0.00%	nouns	16.34%	24.45%	
26	159	_PN	0.44%	0.10%				
27	2585	_PP	5.59%	1.67%				
28	563	_PPAS	3.10%	0.36%				
29	755	_PPAS	0.93%	0.49%				
30	110	_PPL	0.22%	0.07%				
31	81	_PPLS	0.08%	0.05%				
32	113	_PPO	1.44%	0.07%				
33	205	_PPDS	0.41%	0.13%	pronouns	12.22%	2.95%	
34	4024	_VB	4.37%	2.60%				
35	1363	_VBD	4.96%	0.88%				
36	1491	_VBG	1.74%	0.96%				
37	5284	_VBN	2.42%	3.42%				
38	1449	_VBZ	0.27%	0.94%	full verbs	13.75%	8.80%	
39	ABS(J)	Tag	% pro K-Wort	% pro J-Worte		K	J	

• Als Diagramm ausgeben

