



Korpuslinguistik
Annis₂-Korpussuchtool
Suchen in tief annotierten Korpora

SE Historische Korpora, 26.06.2012
Prof. Anke Lüdeling, Carolin Odebrecht
Präsentation: Malte Belz, malte.belz@hu-berlin.de

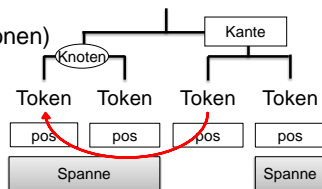
Überblick

- Was ist Annis?
- Wie sieht das Webinterface aus?
- Wie werden Suchanfragen formuliert?

1

Was ist Annis?

- **ANNIS** steht für
 - **ANN**otation of **I**nformation **S**tructure
 - <http://www.sfb632.uni-potsdam.de/d1/annis/>
- Suchmaschine für tief annotierte, multimodale Korpora
 - Token (-annotationen)
 - Spannen
 - Bäume
 - Pointer



Links

- **Annis-Portal (öffentlich)**
<http://korpling.german.hu-berlin.de/Annis/search.html>
 - **Annis-Portal speziell für Lernerkorpora (Falko)**
<http://korpling.german.hu-berlin.de/falko-suche>
 - **Jetzt: Übung mit**
 - **Annis-Portal für diachrone Korpora (DDD)**
<http://korpling.german.hu-berlin.de/ddd/search.html>
- Bitte jetzt diesem Link folgen

Ziele der heutigen Sitzung

<http://korpling.german.hu-berlin.de/ddd/search.html>

- Wie und was kann man in ANNIS suchen?
 - Linguistische Muster
 - Token-Annotationen (Lemmata/Wortarten)
 - Spannenannotationen
 - Syntaktische Annotationen
 - Konstituenten
 - Abhängigkeiten
 - Annis findet nur das, was annotiert ist!!!
- Wie sucht man nach mehreren/beliebigen Annotationen gleichzeitig?
- Wie filtert man nach Metadaten?

4

Das Web-Interface: Tutorial

ANNIS - Tutorial

Search Form

ANNISQL: word="Globe"

Query Builder: Show as

Result: VMS Query

Micro Corpora	Name	Words	Tokens
<input type="checkbox"/>	FalkoEssay1V2_0	95	7000
<input type="checkbox"/>	FalkoEssay2V2_0	248	131599
<input type="checkbox"/>	FalkoSummary1V1_2	57	21211

Search: Export

Context Left: 5

Context Right: 5

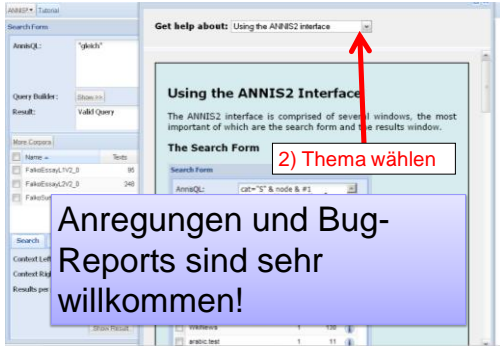
Results per page: 10

Show Result

1) Tutorial öffnen

5

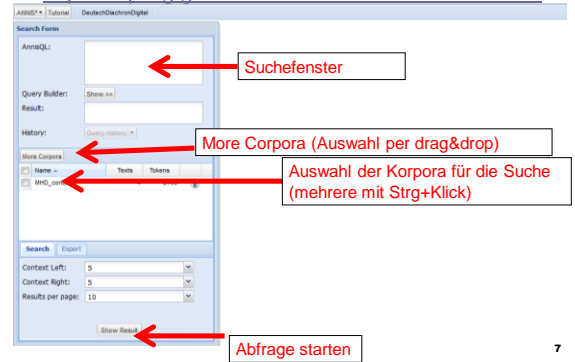
Das Web-Interface: Tutorial



6

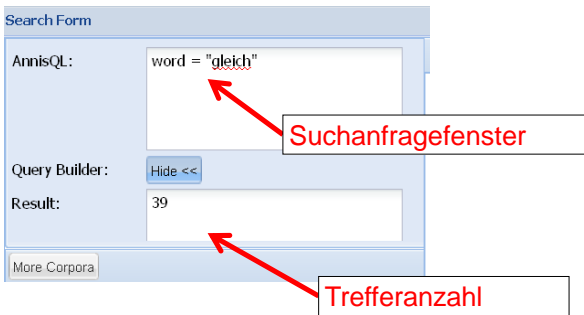
Das Web-Interface: Abfrage

<http://korpling.german.hu-berlin.de/ddd/search.html>



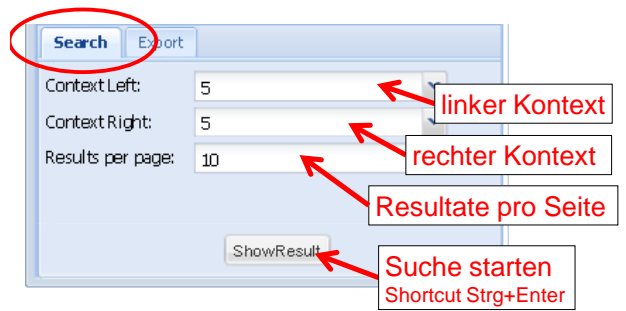
7

Das Web-Interface: Suchfenster



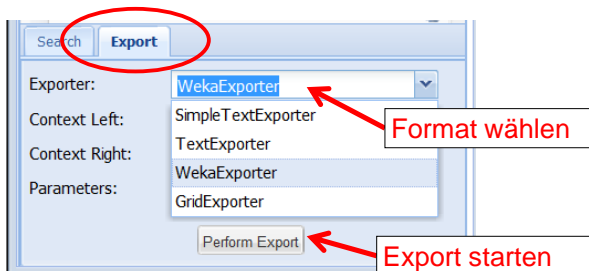
8

Das Web-Interface: Such-Einstellungen



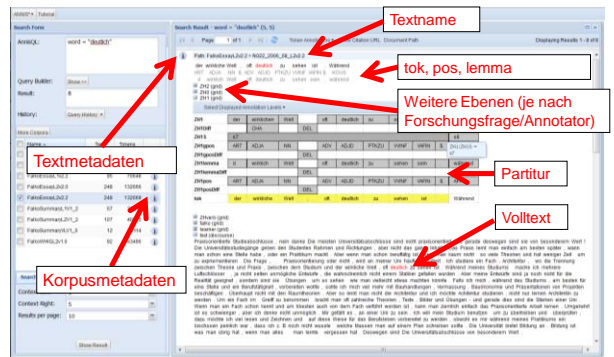
9

Das Web-Interface: Export-Einstellungen



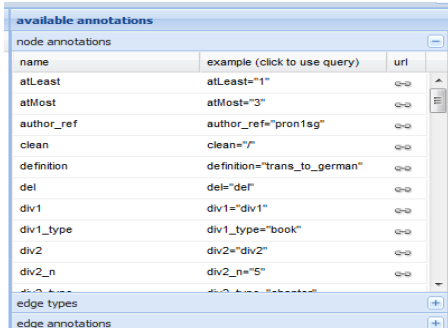
10

Das Web-Interface: Treffer



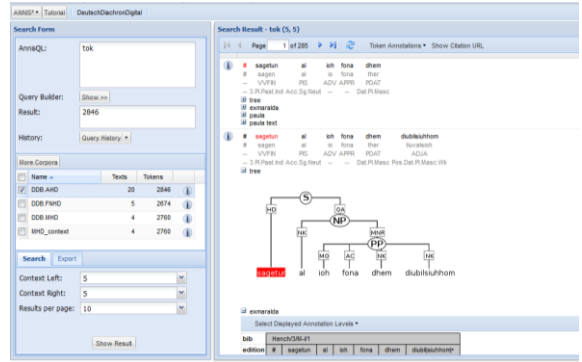
11

Web-Interface: Meta-Information



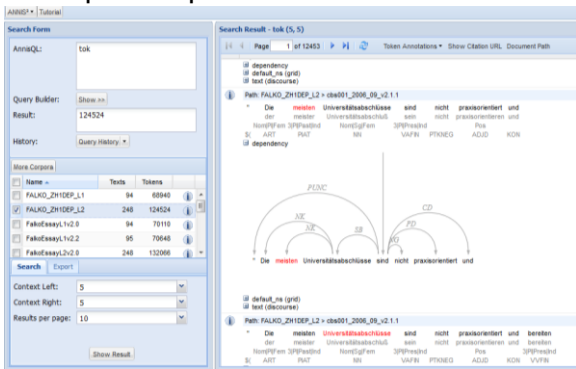
12

Beispiel: Syntaktische Beziehungen



13

Beispiel: Abhängigkeiten



14

Wie suchen wir?

- Prinzip I: Variablen-Wert-Paare (Attribut-Wert-Paare, Layer-Wert-Paare)
- Prinzip II: Relationssuche

15

Prinzip I: Variablen-Wert-Paare

■ tok = "das"

Variable (Layer) Wert

tok	Sofern	das	System	herrscht
pos	KOUS	ART	NN	VVFIN
lemma	sofern	d	System	herrschen

Anfrage in Anführungszeichen: findet exakt diesen String!

16

Prinzip I: Variablen-Wert-Paare

■ pos = "NN"

Variable Wert

tok	Sofern	das	System	herrscht
pos	KOUS	ART	NN	VVFIN
lemma	sofern	d	System	herrschen

17

Prinzip I: Variablen-Wert-Paare

■ lemma = "d" ...findet die, dem, den, ...

Variable3 ("Lemma")

Wert

word	Sofern	das	System	herrscht
pos	KOUS	ART	NN	VVFIN
lemma	sofern	d	System	herrschen

18

beliebig erweiterbar...

■ satz = "NS"

Variable4 ("Satztyp")

Wert

word	Sofern	das	System	herrscht
pos	KOUS	ART	NN	VVFIN
lemma	sofern	d	System	herrschen
satz	NS			

19

beliebig erweiterbar...

■ satz = "NS"

...findet alle Nebensätze wie
Sofern das System herrscht

(sofern die Daten wie gezeigt annotiert sind)

20

Aufgabe 1

<http://korpling.german.hu-berlin.de/ddd/search.html>

- Laden des Korpus „ridges.herbology“ (more corpora → drag&drop → check)
- Suchen Sie nach allen Vorkommen des Tokens „Kraut“

tok = "Kraut"

21

The screenshot shows the ANNO3 search interface. On the left, the search form has 'tok="Kraut"' entered. Below it, a table lists corpora: 'ridges.herbology' with 14 texts and 62724 tokens. The main search results area shows several entries for the token 'Kraut' with their respective grammatical annotations and document views.

22

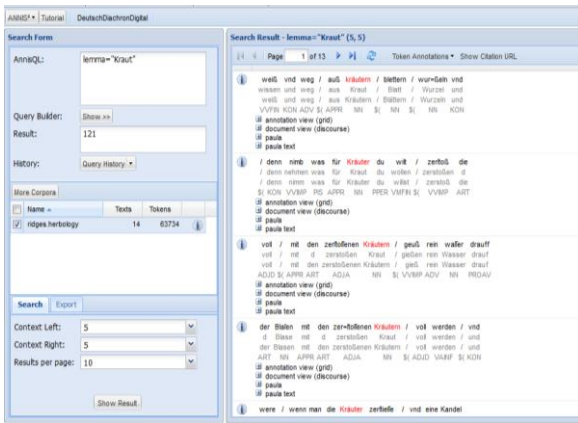
Aufgabe 2

- Suchen Sie nach allen Vorkommen des Lemmas „Kraut“

lemma = "Kraut"

□ Ergebnis: 121

23



24

Mustersuche (reguläre Ausdrücke)

- Annis₂ erlaubt Mustersuchen auf allen Annotationsebenen
- Mustersuchen werden statt in " " in // eingefügt
- Z. B. kann man damit nach allen Wörtern suchen, die „kraut“ enthalten.

```
tok = /. *kraut.*/
```

25

Operatoren: Joker .

- ein beliebiges Zeichen al. → *als, alt, ...*
- ■ zwei beliebige Zeichen al.. → *alle, alte, also*
- ■ ■ drei beliebige Zeichen al... → *alles, altes, alias, ...*

Aufgabe 4

- Welche Wortformen erhalten Sie?

```
tok = /g.b./
```

Lösung: *gibt, gebe* (in ridges.herbology).
Denkbar wäre auch *gäbe*

27

Operatoren: ? und * +

- das[?] das vorherige Zeichen ist optional
→ ϕ, s → *da, das*
- das^{*} das vorh. Zeichen kommt 0- bis ∞ mal vor
→ ϕ, s, ss, \dots → *da, das, dass, dasssss*
- das⁺ das vorh. Zeichen kommt 1- bis ∞ mal vor
→ *s, ss, \dots* → *das, dass, dasssss*

28

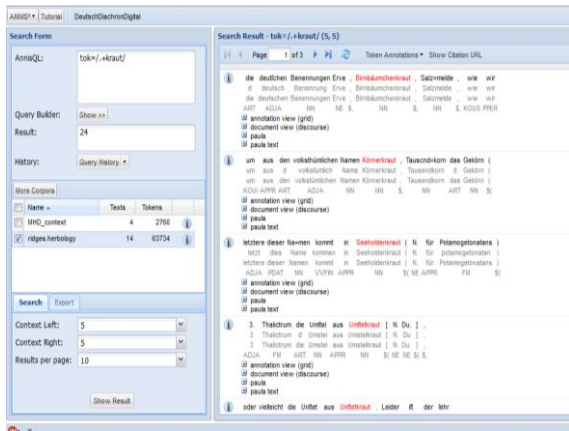
Beispiel 5

- Versuchen Sie alle Token zu finden, die auf *kraut* enden.

```
tok = /.+kraut/
```

Findet *Birnbäumchenkraut, Körnerkraut, Wielandskraut...*

29



30

Beispiel 6

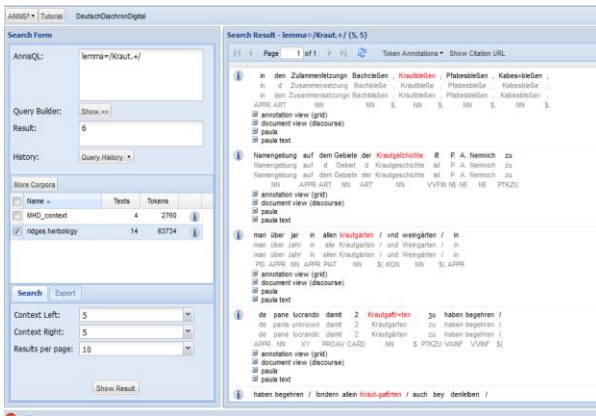
- Versuchen Sie alle Lemmata zu finden, die mit *Kraut* beginnen.

lemma = /Kraut.+/

Besonderheiten der Lemma-Ebene:

- *Kraut* wird zu Beginn großgeschrieben (keine Alternativensuche notwendig)
- Es kann passieren, dass der Tagger das Lemma nicht kennt und mit <unknown> annotiert hat, falls diese Ebene nicht manuell korrigiert wurde

31



32

Alternativen: a oder b = (a|b)

- Mit Klammern und | ("oder") kann man gleichzeitig nach verschiedenen Wörtern suchen:

tok = /(Mann|Frau|Kind)/

Was ist mit *better*?
Suche in der clean- oder norm-Ebene!

- Nach verschiedenen Formen:

tok = /(Mann|Mannes)/

- Oder Zeichenketten:

tok = /bes(ser)t.?!/

besser, bessere
best, beste
Und andere Strings:
bests...

33

Aufgabe 7

- Finden Sie alle Formen des Nomens *Kraut* im Plural
- Welche Formen können in älteren Texten vorkommen?
 - Kräuter, kräuter, Kräutter, krütter, Kreuter, kreuter, Kreutter, kreutter, Kräutern, Kreutern, kreüter, Kräuteren...
 - (ohne Komposita)

34

Lösung 7

Nicht im aktuellen ridges v.1

tok=/(K|k)rr?(eu|äü|ä|äü|eü|äü)tt?e?r?e?n?./

oder

tok=/(K|k)rr?(e(u|ü)|ä(ü|ä)|ä(u|ü)tt?e?r?e?n?./

→Häufig gibt es alternative Suchanfragen für dieselben Treffermengen.
→Unbekannte Schreibweisen können nicht antizipiert werden! (*krühter*)

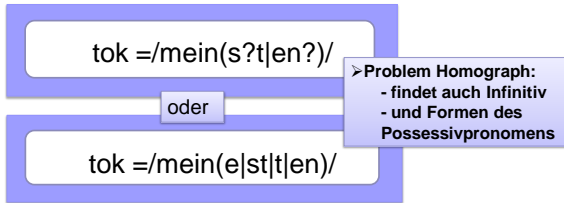
Einfachere Suche auf Norm oder Lemma → Export

35

Beispiel 8

mein	e
mein	st
mein	t
mein	en
mein	t
mein	en

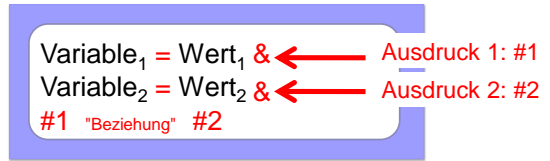
- Man möchte alle konjugierten Formen des Verbs *meinen* im Präsens finden, aber keine anderen Formen.



36

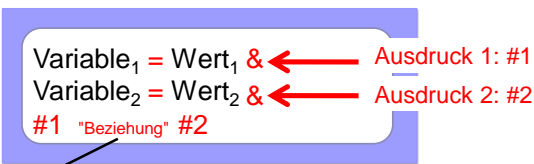
Prinzip II: Relationen

- Einzelne Variable-Wert-Paare werden durch "&" verbunden.
- Zwischen den Paaren muss **IMMER eine Beziehung hergestellt** werden
- Auf die VW-Paare bezieht man sich mit # der Reihe nach.



37

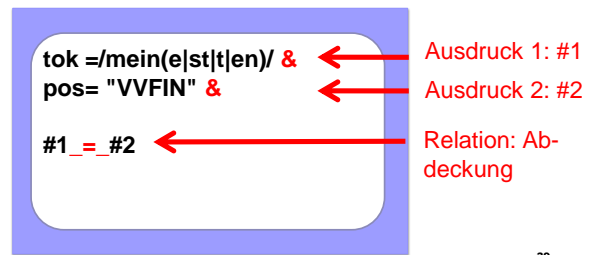
Prinzip II: Relationstypen



Operator	Description	Illustration	Notes
.	direct precedence	A B	For non-terminal nodes, precedence is determined by the right most and left most terminal children
.*	indirect precedence	A x y z B	For specific sizes of precedence spans, .n.a can be used, e.g. .3,4 - between 3 and 4 tokens
⊆	identical coverage	A B	Applies when two annotation cover the exact same span of tokens
⊇	inclusion	A B	Applies when one annotation covers a span identical to or larger than another
>	direct dominance	A 1 B	A specific edge type may be specified, e.g.: >acedge to find secondary edges. Edges labels are specified in brackets, e.g.: >(ance="a") for an edge with the

Zu Beispiel 8

- Finden Sie nun Vorkommen von tok =/mein(e|st|t|en)/, die ausschließlich finite Vollverben sind.



39

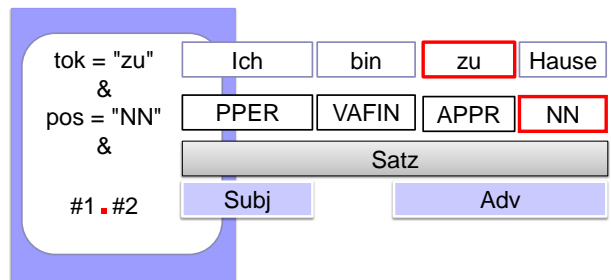
Negation !=

- ! bedeutet Negation
 - Der Operator wird vor dem "="-Zeichen eingefügt.
 - Finden Sie in alle Vorkommen von tok =/mein(e|st|t|en)/, die nicht als finites Vollverb getaggt sind



40

Suche nach Abfolgen:
z.B. Nomen folgt auf "zu"



41

Tokenfolgen – Aufgabe 9

- Suchen Sie nach zwei aufeinanderfolgenden Adjektiven.
- Achtung: Es gibt zwei Typen von Adjektiven
 - ADJA & ADJD

```
pos = /ADJ./ &
pos = /ADJ./ &
#1.#2
```

42

Metadaten

- Sie können Texte nach Metadaten filtern. Metadaten finden Sie unter dem i-Button in den Suchergebnissen
- Metadaten müssen nur verknüpft, aber nicht relativiert werden
- Beispiel: alle Token von Texten, die in Frankfurt publiziert wurden

```
tok &
meta::pubPlace="Frankfurt"
```

43

Tokenfolgen und Metadaten

- Suchen Sie nach zwei aufeinanderfolgenden Adjektiven in einem Text von 1603!

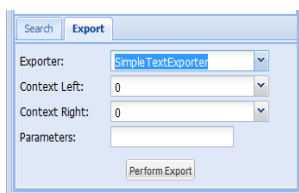
```
pos = /ADJ./ &
pos = /ADJ./ &
#1.#2 &
meta::date="1603"
```

44

Beispiel Export

- Sie suchen alle Schreibvarianten von *dass*
- Dazu können sie
 - nach lemma= "dass" suchen
 - die gefundenen Ergebnisse zur besseren Betrachtung aller Vorkommen exportieren

45



Ergebnisse (Auszug):

1. [daf]	21. [daß]	80. [dafs]
2. [dalf]	22. [dass]	81. [dafs]
3. [dalf]	23. [daß]	82. [dafs]
4. [dalf]	24. [daß]	83. [dafs]
5. [dalf]	25. [daß]	84. [dafs]
6. [dz]	26. [daß]	85. [daß]
7. [daf]	27. [dass]	319. [Das]
8. [daf]	28. [dass]	320. [dafz]
9. [daf]	29. [dass]	321. [daß]
10. [Das]	30. [daß]	322. [das]

46

Zusammenfassung Operatoren

- `.` Ein beliebiges Zeichen
- `?` 0 oder 1 Zeichen (des vorherigen Elementes)
- `*` 0 bis unendlich viele Zeichen (d. vorh. E.)
- `+` 1 bis unendlich viele Zeichen (d. vorh. E.)
- `\\` wörtlich (folgendes Zeichen)
- `!` nicht
- `[abc]` Menge (oder `[^abc]=alles außer abc`)
- `(a|b)` a oder b (auch: `[ab]`)
- `a{2,3}` a 2 bis 3mal

49

Zusammenfassung Relationen

Vielen Dank!

- **&** verbindet Suchanfragen
- **#1.#2** direkte Präzedenz
- **#1.*#2** indirekte Präzedenz
- **#13,5#2** #2 folgt mit 3 bis 5 Einheiten Abstand
- **#1_=#2** #1 und #2 deckungsgleich
- **#1_o_#2** #1 und #2 überlappen sich
- **#1_i_#2** #1 inkludiert #2
- Weitere Informationen

□ <http://www.sfb632.uni-potsdam.de/d1/annis/>

50

51