

ANNIS Einführung RIDGES Herbology 4.1

Ridges
Register in Diachronic
German Science

Laura Perlitz und Carolin Odebrecht Humboldt-Universität zu Berlin



Korpus RIDGES Herbology



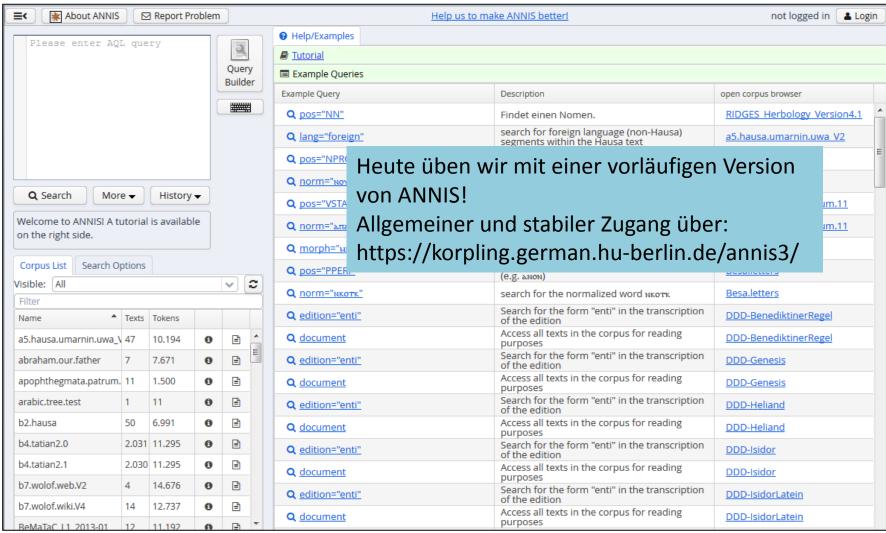
- Wiederholung
 - Kräuterkundekorpus
 - 15.-19. Jahrhundert
 - verschiedene Dialekte aus dem deutschsprachigen Raum
 - Arbeit von drei Seminaren aus B.A. und M.A. von ca. 50 Studenten
 - Projekthomepage
 - http://korpling.german.huberlin.de/ridges/documentation v4.1 en.html
 - Annotationsrichtlinie
 - http://korpling.german.huberlin.de/ridges/download/pubs/annotationGuidelines v4.1
 .pdf

ANNIS Such- und Visualisierungstool



- webbasiertes Suchtool f
 ür Korpora
 - Projekthomepage http://annis-tools.org/
 - Zugang auf das ANNIS-Tool https://korpling.german.hu-berlin.de/annis3
- Was kann ANNIS?
 - findet Annotationen in einem Korpus
 - > findet mehrere Annotationen in bestimmter Relation zu einander
 - > gibt die gesuchten Annotationen als Treffer in vielfältigen Visualisierungen aus
 - > exportiert diese Treffer
 - > u.v.m.
- ➤ Hier im Seminar lernen wir mit dem Korpus RIDGES Version 4.1, wie man ANNIS benutzt!

https://korpling.german.hu-berlin.de/annis3-snapshot/



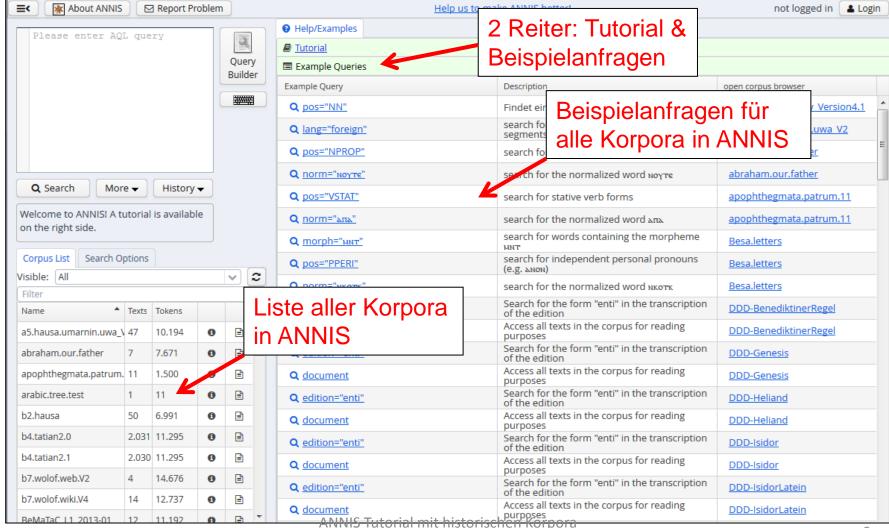


Interface

Startseite, Hilfestellung, Query Fenster

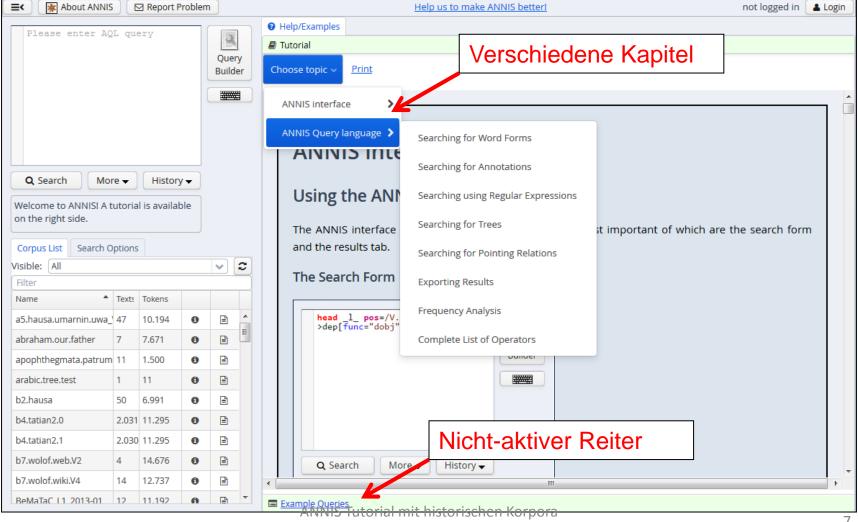
Interface Startseite





Interface **Tutorial**

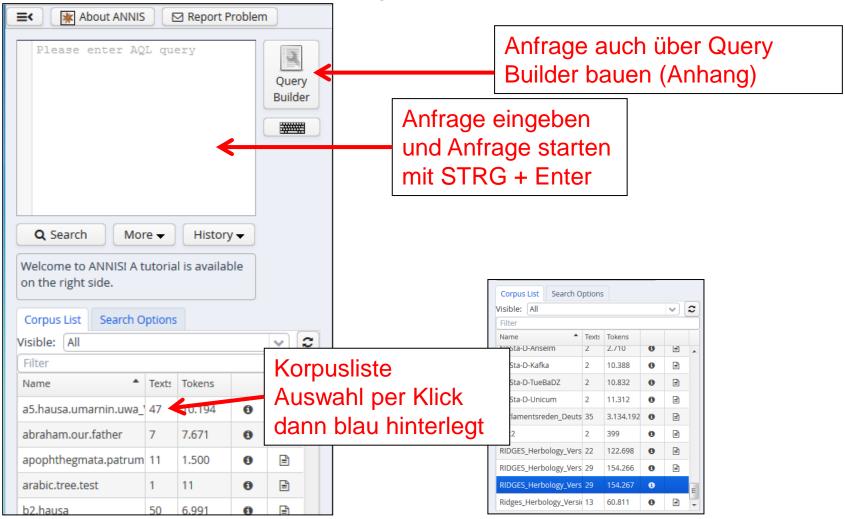




(RIDGES Herbology)

Interface Query-Fenster





ANNIS Query Language



- oder: Wie sage ich ANNIS, was ich suche?
- für die Suche nach Annotationen gelten zwei Prinzipien:
 - Variable-Wert-Paar (VW-Paar)
 - Relationen zwischen VW-Paaren

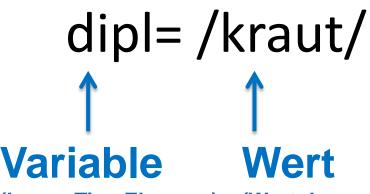
Token

- 1) Als Token bezeichnet man häufig die kleinste (technische) Einheit in einem Korpus!
- 2) Ein Token entspricht oft (aber nicht immer) einem orthographischen Wort oder Satzzeichen!
- Nach diesen Einheiten kann man in ANNIS suchen.

Token Token Token Token Token ...

dipl	wider	auff	ein	friſches	zerítolíenes	kraut	1	thue	es	wider	in
clean	wider	auff	ein	frisches	zerstossenes	kraut	1	thue	es	wider	in
norm	wieder	auf	ein	frisches	zerstoßenes	Kraut	1	tue	es	wieder	in

Prinzip I: Variable-Wert-Paar



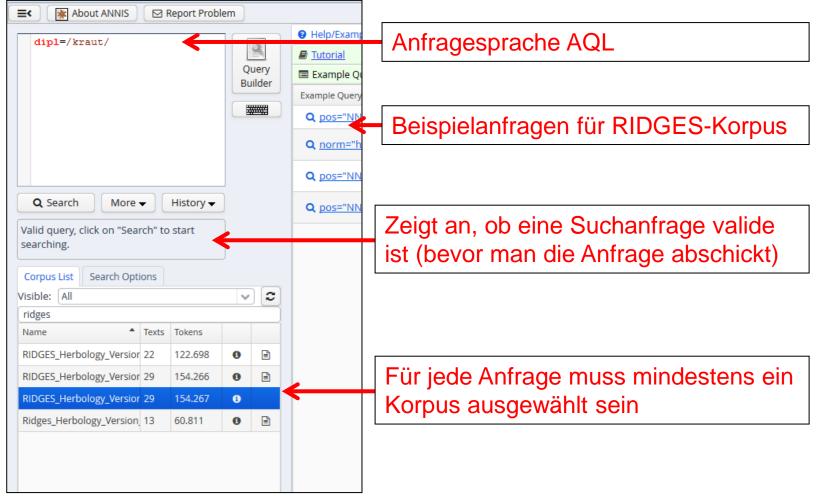
- Voraussetzung ist das Vorhandensein einer Ebene namens "dipl". (Metadaten!)
- 2) Erwartetes Ergebnis ist es, exakt alle Vorkommen dieser Zeichenkette in "dipl" im ausgewählten Korpus zu finden.

(Layer, Tier, Ebene ...) (Wort, Lemma, Satz, Wortart ...)

dipl	wider	auff	ein	friſches	zerítollenes	kraut	1	thue	es	wider	in
clean	wider	auff	ein	frisches	zerstossenes	kraut	1	thue	es	wider	in
norm	wieder	auf	ein	frisches	zerstoßenes	Kraut	1	tue	es	wieder	in

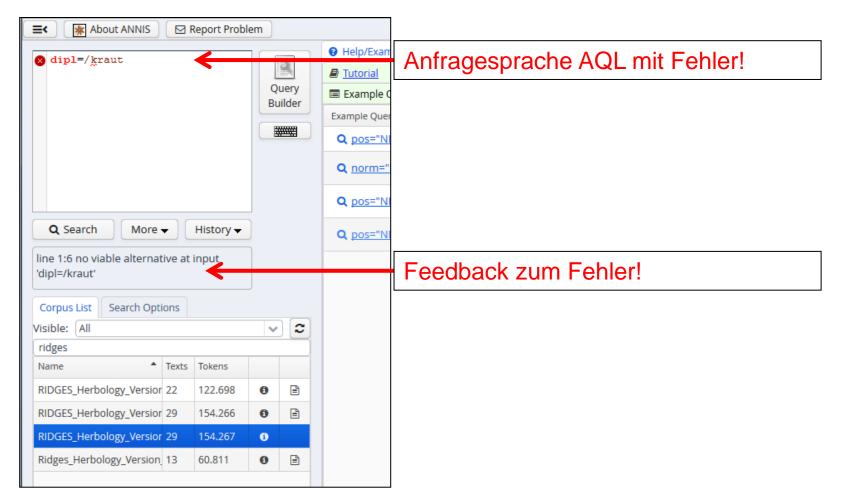


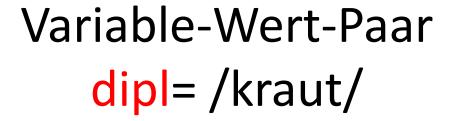




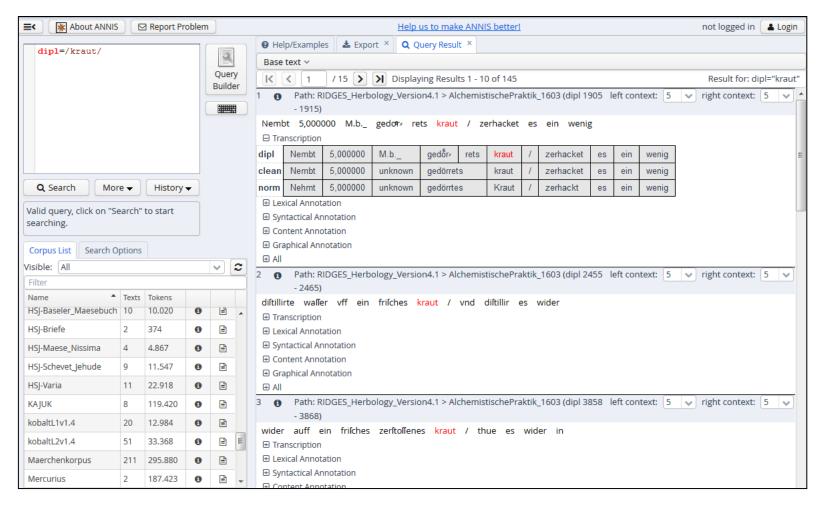












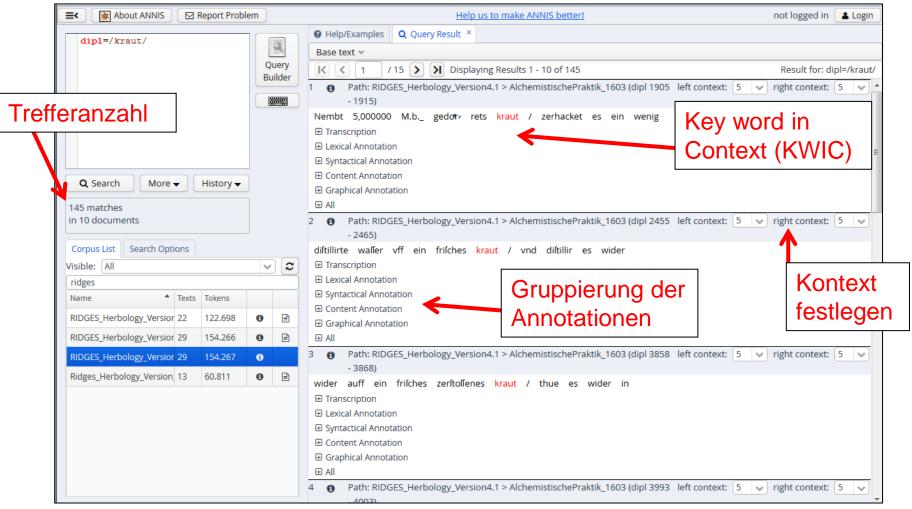


Interface

Treffer, Visualisierung, Metadaten, Suchverlauf

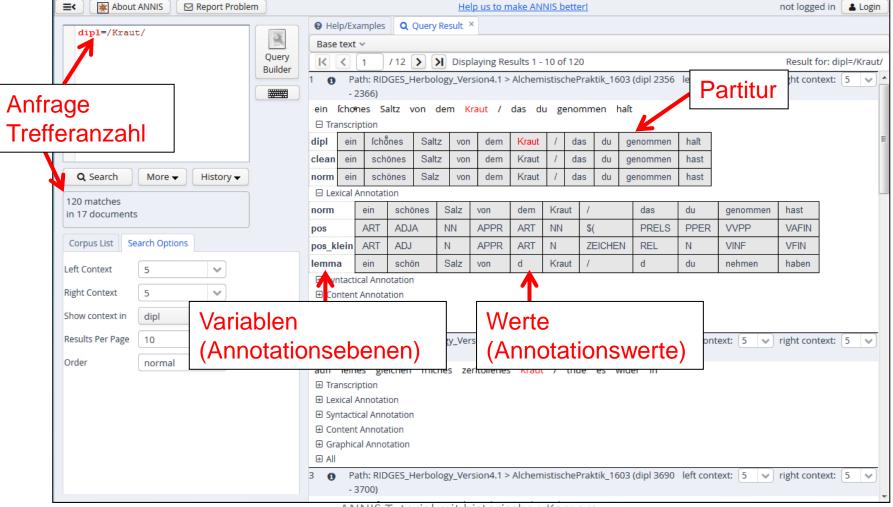
Interface Treffer







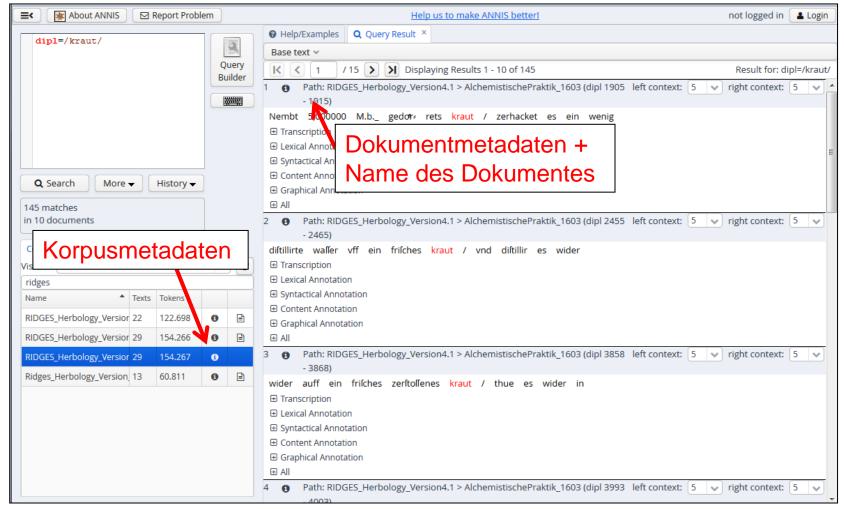




ANNIS Tutorial mit historischen Korpora (RIDGES Herbology)

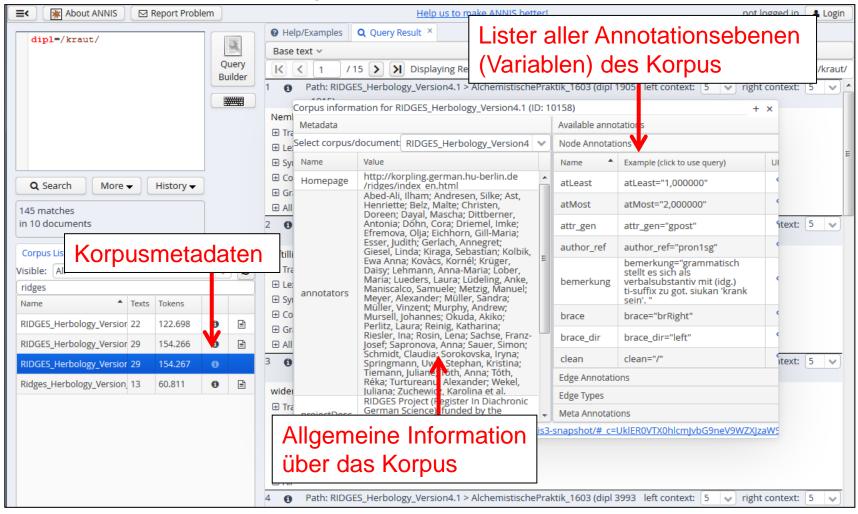






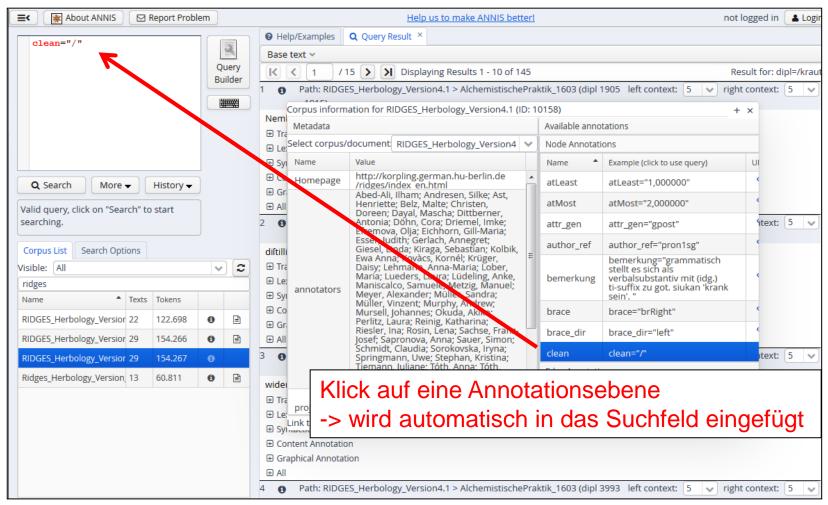
Interface Korpusmetadaten





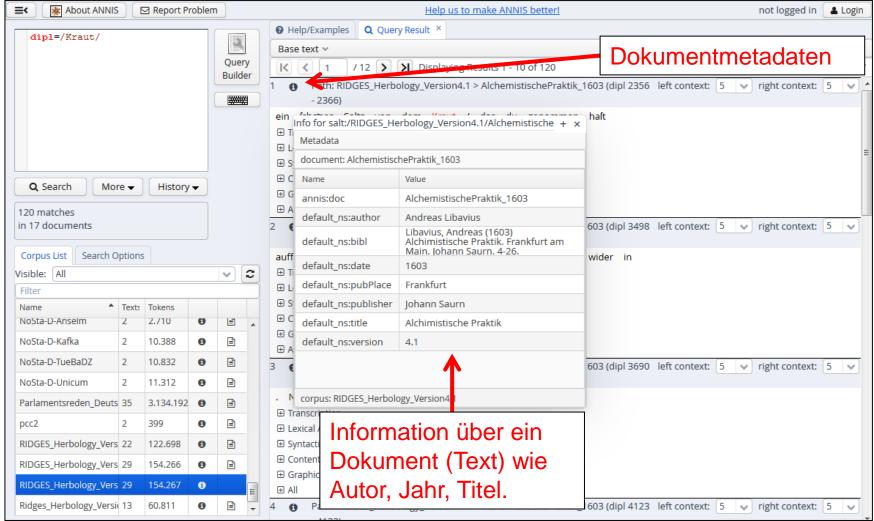






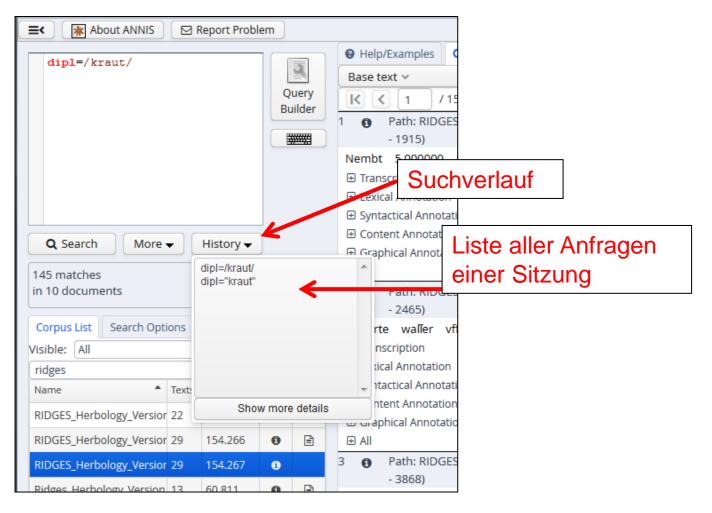
Interface Dokumentmetadaten

















Prinzip I: Variable-Wert-Paar

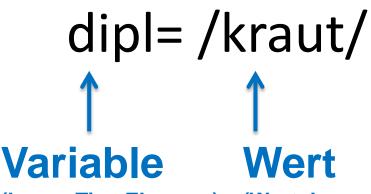
- AQL
 - Namen der Annotationsebene → Variable dipl =
 - annotierte Kategorien

→Wert

/kraut/

- Wissen, welche Annotationen vorhanden sind
 - Korpusmetadaten
 - Annotationsrichtlinien des Korpus (für RIDGES vorhanden http://korpling.german.hu-berlin.de/ridges/download/pubs/annotationGuidelinesv4.1.pdf)

Prinzip I: Variable-Wert-Paar



- 1) Voraussetzung ist das Vorhandensein einer Ebene namens "dipl".
- 2) Erwartetes Ergebnis ist exakt alle Vorkommen dieser Zeichenkette in "dipl" im gesamten Korpus zu finden.

(Layer, Tier, Ebene ...) (Wort, Lemma, Satz, Wortart ...)

Y											
dipl	wider	auff	ein	friſches	zerítollenes	kraut	/	thue	es	wider	in
clean	wider	auff	ein	frisches	zerstossenes	kraut	/	thue	es	wider	in
norm	wieder	auf	ein	frisches	zerstoßenes	Kraut	1	tue	es	wieder	in

Historische Texte



- Varianz der Orthographie bzw. Setzung
- vieles nicht wirklich "vorhersehbar"
- graphische/(text-) strukturelle Informationen
 - Zeilenumbruch,
 Seitenumbruch,
 Setzung

Rräuter	Kräuter	Deutsche Pflanzennamen 1870
Rrater	Kreuter	Alchmistische Praktik 1603
Kräuteren	Krauteren	Pflantz.Gart. 1639
Freüter	kreüter	Fuchs New Kreüterbuch 1543
Rreutter	Kreutter	Alchimistische Praktik
Kreutern	Kreutern	Alchimistische Praktik 1603
Fraut	kraut	Alchimistische Praktik 1603
Kraut	Kraut	Alchimistische Praktik 1603
Rräutern	Krautern	Alchimistische Praktik 1603

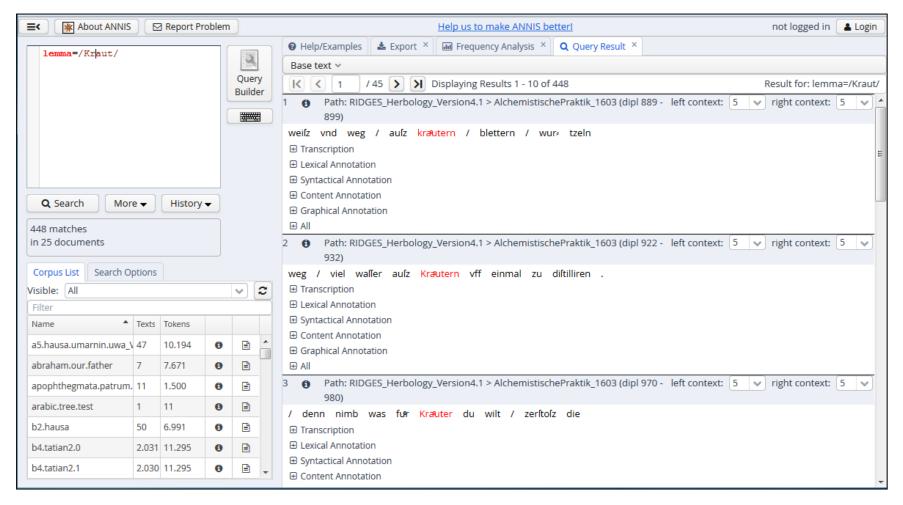
Aufgabe Schreibvarianten



- Finden Sie alle Schreib- und Flexionsvarianten von Kraut!
 - für die Suche nach allen Varianten wählen Sie die passende Annotationsebene (Variable) aus
 - passende Variable wäre hier lemma
 - und setzen den Wert /Kraut/ ein
 - ➤lemma=/Kraut/











 Schauen Sie sich die Treffer zu lemma=/Kraut/ genau an!

weiſz	vnd v	veg /	auſz	kræ	utern /	blettern	/	wur៖ tzeli	n		
☐ Transcription											
dipl	weiſz	vnd	weg	1	auſz	krautern	1	blettern	/	WUΓ۶	tzeln
clean	weisz	vnd	weg	/	ausz	kräutern	1	blettern	1	wurtze	eln
norm	weiß	und	weg	1	aus	Kräutern	1	Blättern	1	Wurze	eln

– Es werden alle historischen Wortformen von Kraut gefunden!

Aufgabe Normalisierung



- Suchen Sie das Lemma von zusammensetzen!
 - Finden von ganz unterschiedlichen historischen Schreibvarianten!
 - lemma=/zusammensetzen/
 - Treffer 1

```
erîten grade / vnd ſey zůſammen geſetzt auſz widerwertiger ſubſtantz . das
```

Treffer 2

```
, die aus drei Worten zulammen geletzt lind , wie die oben
```

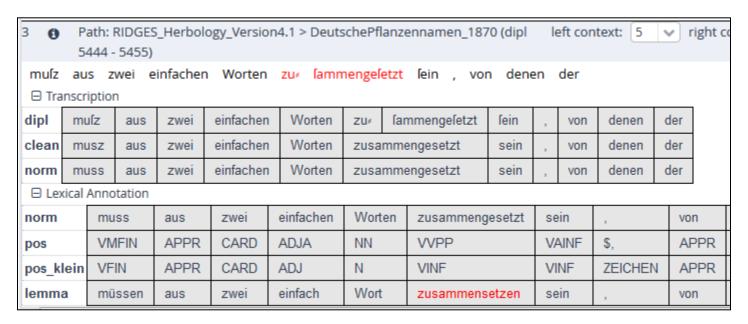
Treffer 3

```
muſz aus zwei einfachen Worten zus ſammengeſetzt ſein , von denen der
```

Aufgabe Normalisierung



- Schrittweise Normalisierung (dipl>clean>norm) erlaubt
 - Zusammenführung von historischen Schreibweisen (dipl)
 - Zuweisung von Wortarten und Lemmatisierungen auf der normierten Ebene (norm)







- 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
- (a|b) a oder b (auch: [ab])





- Welche Ergebnisse erwarten Sie für folgende Anfragen, sogenannte Mustersuchen?
 - norm=/g.b./
 - ➤ gebe, gibt (für RIDGES, theoretisch noch andere möglich)
 - dipl=/r(a|o)t/
 - *>rot, rat* (für RIDGES)
 - dipl=/meint?/
 - > mein, meint (für RIDGES)

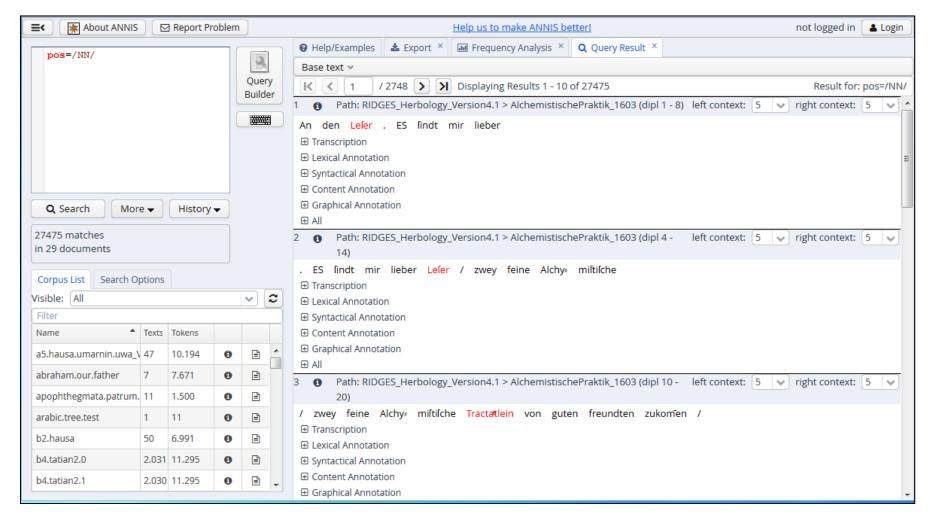
Aufgabe Wortart



- Wie finden Sie alle Appellativa in RIDGES?
 - passende Variable wäre hierpos
 - >pos=/NN/







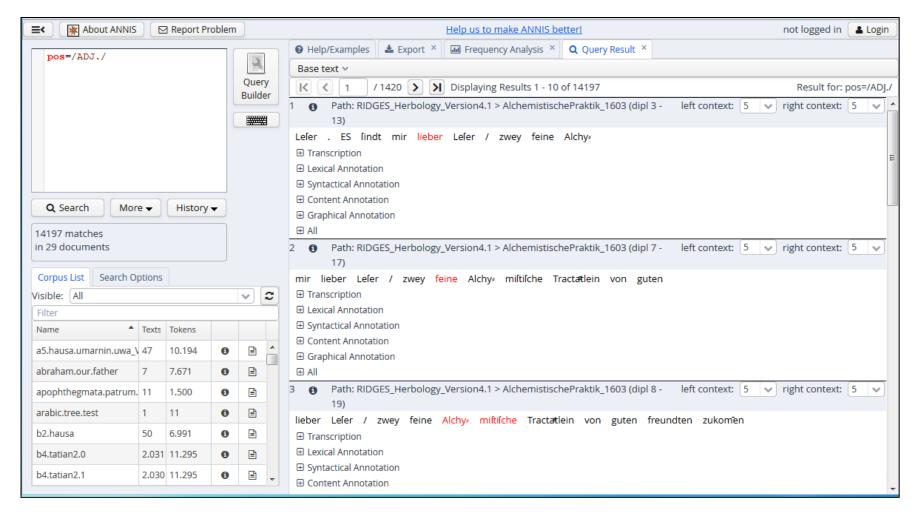
Aufgabe Wortarten



- Wie finden Sie alle Adjektive unabhängig von ihren Bezugswörtern in RIDGES?
 - Was sagen die Richtlinien (STTS)?
 - passende Variable wäre hier pos
 - STTS: Unterschied ADJA und ADJD?
 - ➤ pos=/ADJ./







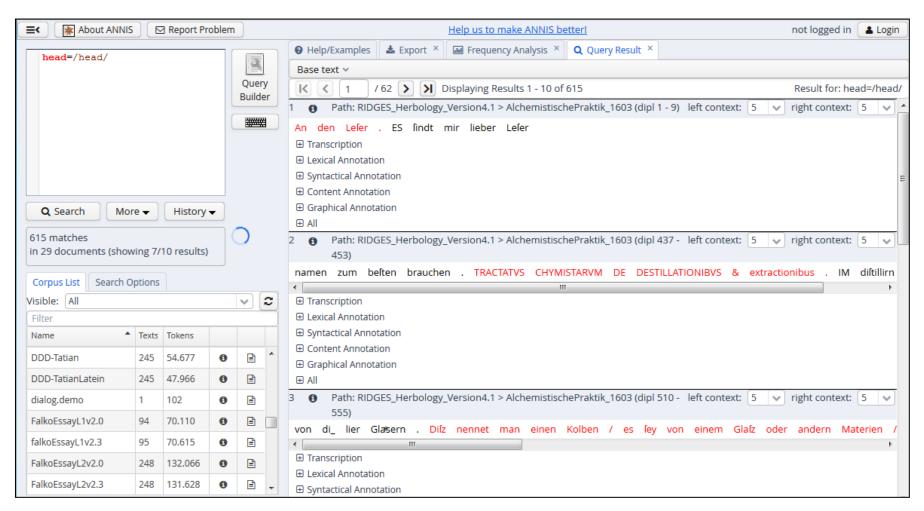
Aufgabe Suchen von Annotationsebenen



- Finden Sie heraus, ob es in RIDGES Herbology
 Annotationen für Überschriften gibt?
 - Wie sieht eine Anfrage aus, um Überschriften zu finden?
 - passende Variable wäre hier head
 - ▶head=/head/

Aufgabe Suchen von Annotationsebenen





Aufgabe Kombinierte Suche

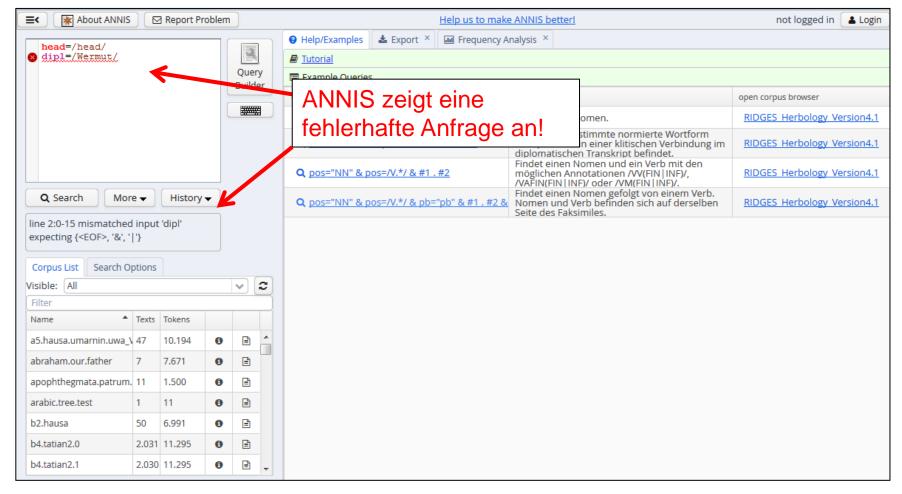


- Finden Sie die diplomatische Wortform
 Wermut, die in einer Überschrift vorkommt!
 - passende Variablen wären hier head und dipl
 - ➤head=/head/
 - ➤ dipl=/Wermut/

Was passiert?







Prinzip II: Relationen

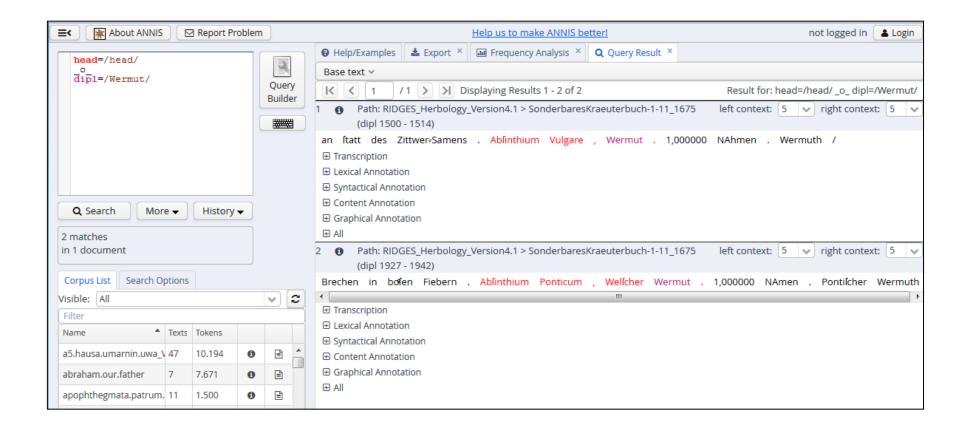
head=/head/ VW-Paar1
o Relation
dipl=/Wermut/ VW-Paar2

- 1) Es gibt mehrere Arten von Annotationen!
- Wissen, wie diese
 Annotationen in Beziehung zu einander stehen können!
- 3) Spannen (head) überlappen Tokenannotationen (dipl)!

XVIII.	Saturn	iinifche Krauter	. Von	ħ	Kratute	rn	ihrer Natu	r				
☐ Transcription												
dipl	XVIII.	Saturninische	Kråuter		Von	ħ	Kråutern	ihrer	Natur			
clean	XVIII.	Saturninische	Kräuter		Von	ħ	Kräutern	ihrer	Natur			
norm	XVIII.	Saturninische	Kräuter		Von	?	Kräutern	ihrer	Natur			
head	head		'									







Syntax-Highlighting



- pro VW-Paar eine Farbe hier unser Beispiel:
 - Farbe Rot für alle Werte der Variable head
 - Farbe Lila für alle Werte der Variable dipl

```
head=/head/
o_dipl=/Wermut/
```

- Treffer erhalten genau diese Farben:
 - mehrere Token rot, da sie zusammen in einer Überschrift stehen
 - ein Token lila, da dieses der gesuchte dipl-Wert

```
Ablinthium Vulgare , Wermut .
```

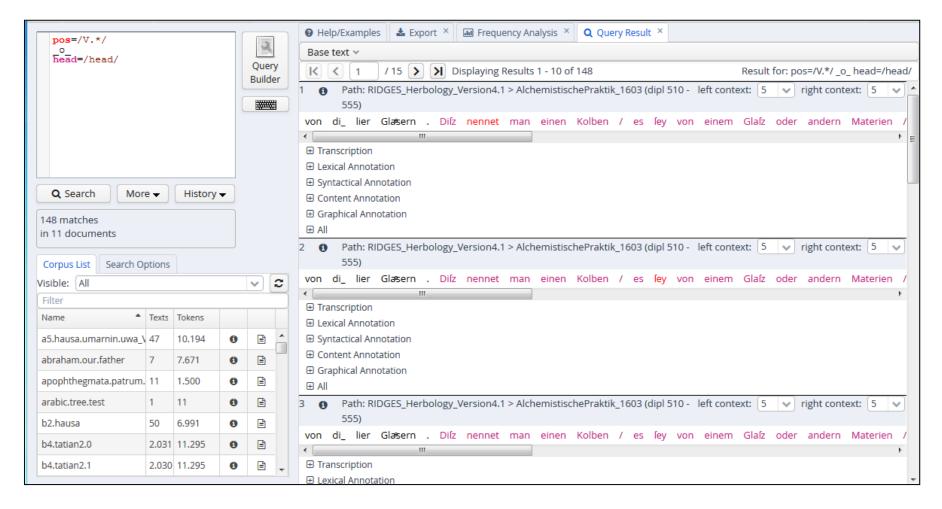
Aufgabe Überlappung



- Suchen Sie ein Verb, das in einer Überschrift vorkommt!
 - passende Variablen wären pos und head
 - Operator _o_
 - > pos=/V.*/ _o_ head=/head/







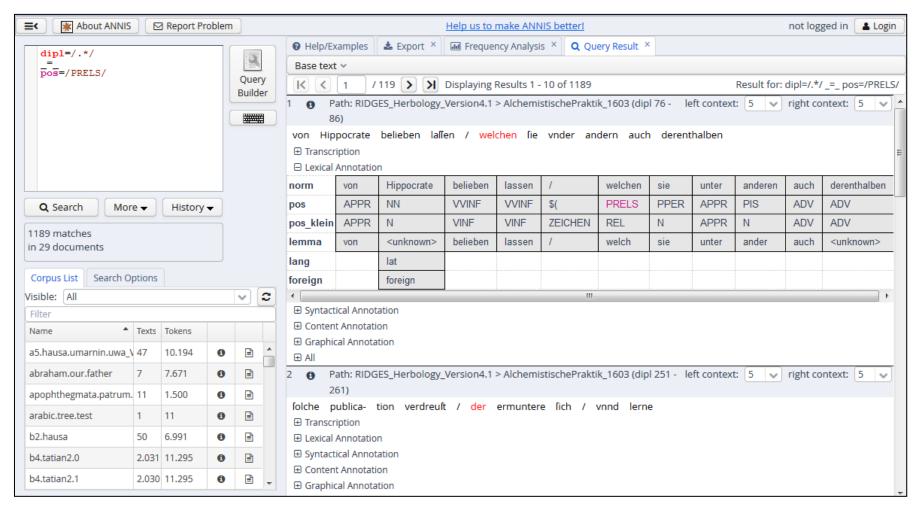
Aufgabe Identität



- Suchen Sie alle diplomatischen Wortformen, die als substituierende Relativpronomen annotiert worden sind!
 - passende Variablen wären dipl und pos
 - Operator _=_
 - >dipl=/.*/ _=_ pos=/PRELS/







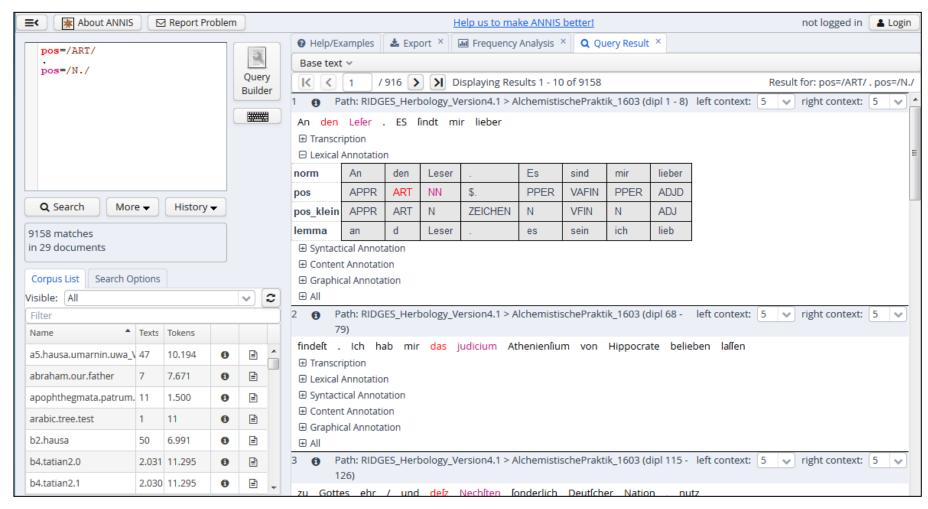
Aufgabe Direkte Präzedenz



- Suchen Sie einen Artikel, der ein Nomen direkt präzediert!
 - passende Variable wäre pos
 - Operator .
 - ➤ pos=/ART/ . pos=/N./







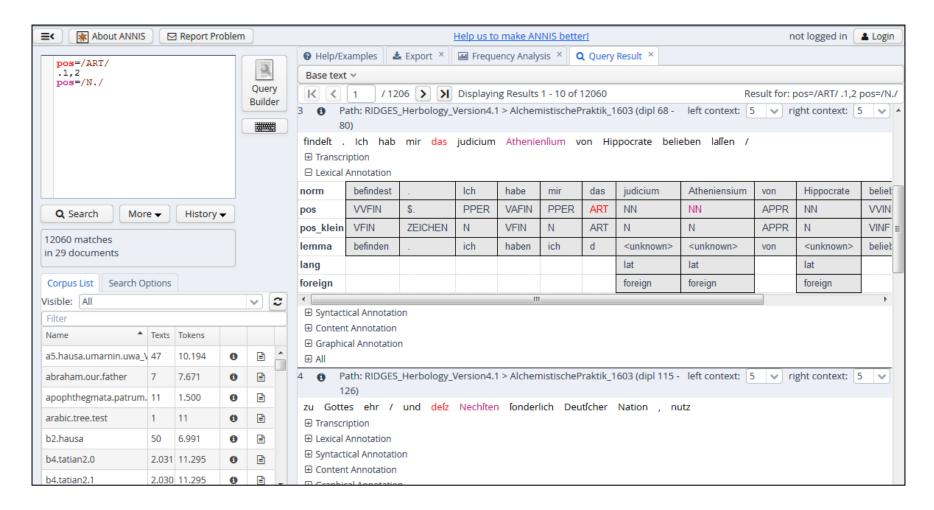
Aufgabe Indirekte Präzedenz



- Suchen Sie einen Artikel (A), der ein Nomen (B) indirekt präzediert. Sie wollen ebenfalls einen möglichen pränominalen Modifikator (C) in ihrer Trefferliste abfangen!
 - passende Variablen wäre pos
 - Operator . und Abstand 1,2 (zu lesen: Ich suche A und B direkte aufeinanderfolgend oder es kann ein Token C dazwischen stehen)
 - >pos=/ART/ .1,2 pos=/N./







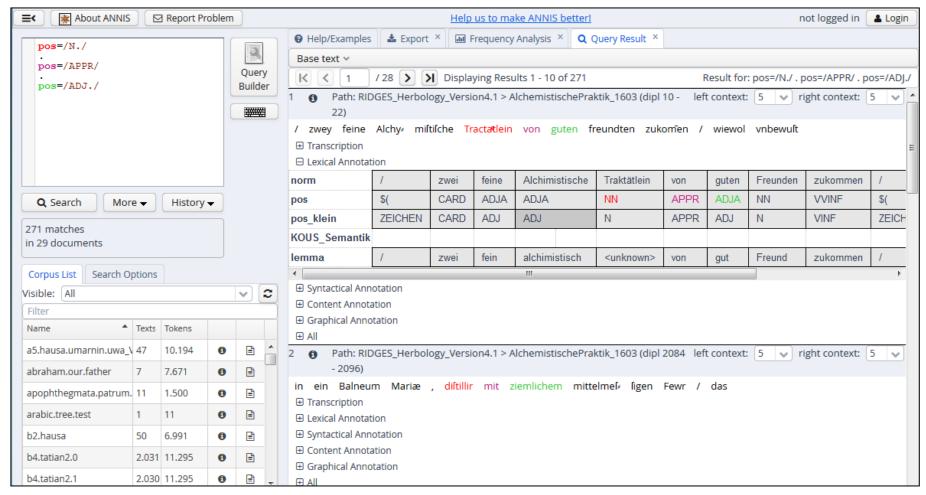
Aufgabe Abfolgen



- Suchen Sie ein Nomen direkt gefolgt von einer Präposition direkt gefolgt von einem Adjektiv!
 - passende Variable wäre pos
 - Operator .
 - >pos=/N./.pos=/APPR/.pos=/ADJ./







Aufgabe Abfolgen



- Suchen Sie ein Nomen direkt gefolgt von einer Präposition direkt gefolgt von einem Adjektiv!
 - passende Variable wäre pos
 - Operator .

/ zwey feine Alchy₅ miſtiſche Tractatlein von guten freundten zukomen / wiewol vnbewuſt

⊞ Transcription

□ Lexical Annotation

norm	1	zwei	feine	Alchimistische	Traktätlein	von	guten	Freunden	zukommen
pos	\$(CARD	ADJA	ADJA	NN	APPR	ADJA	NN	VVINF

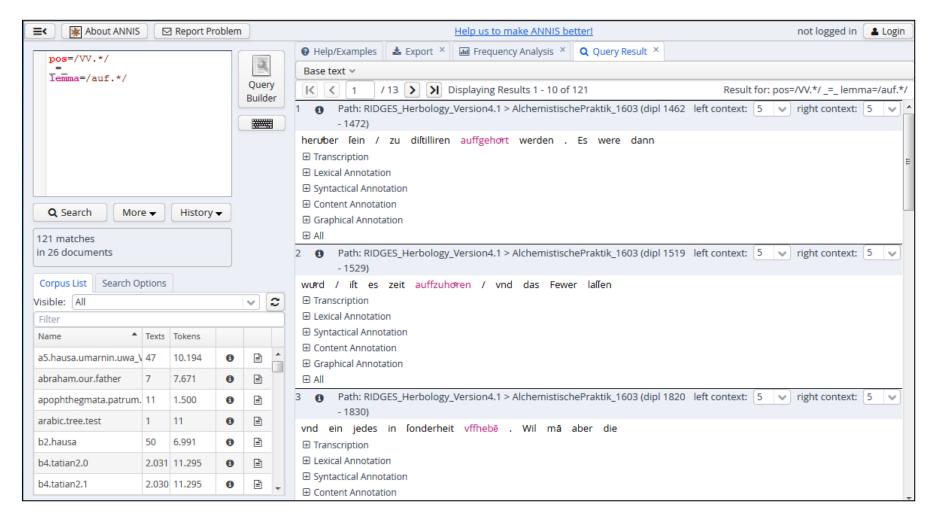
Aufgabe Identität



- Suchen Sie Partikelverben, die die Partikel auf beinhalten!
 - passende Variablen wären lemma und pos
 - Operator _=_
 - >pos=/VV.*/ _=_ lemma=/auf.*/







Aufgabe Identität



- Suchen Sie Partikelverben, die die Partikel auf beinhalten!
 - passend wären lemma und pos
 - Operator _=_
 - >pos=/V.*/_=_ lemma=/auf.*/

Was findet man damit nicht?

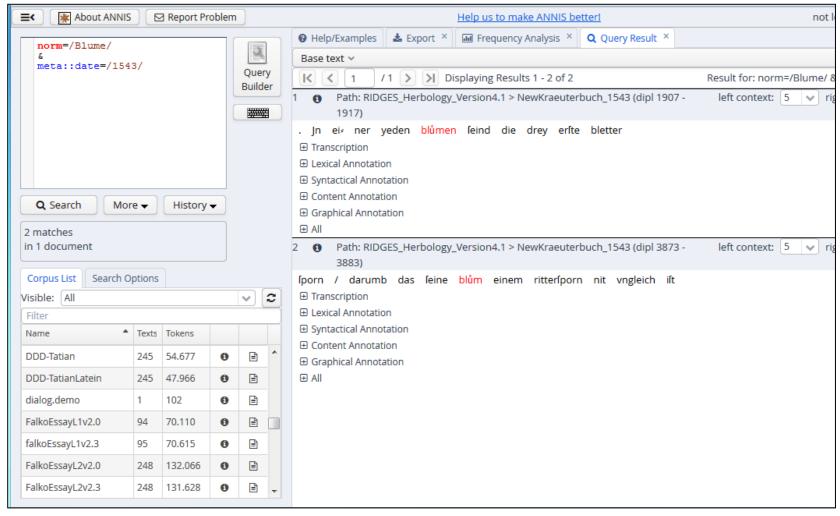
Aufgabe Metadaten



- Finden Sie heraus, ob die normierte Wortform Blume in einem Text aus dem Jahr 1543 zu finden ist!
- Wenn ja, wie oft?
 - passende Variablen wären norm und meta::date
 - Verknüpfungszeichen &
 - ➤ norm=/Blume/ & meta::date=/1543/
 - **≥**2 Treffer







Tipps



- Die Trefferliste zeigt im Kontext nur die Annotationsebenen (Variablen) an, die auch genau in diesem Trefferkontext annotiert worden sind!
- Schauen Sie in die Korpusmetadaten und Annotationsrichtlinien, um zu wissen, welche Annotationsebenen im Korpus vorhanden sind!
- Nicht alle Dokumente in einem Korpus müssen die gleichen (Anzahl und Typ) Annotationsebenen besitzen!

Zusammenfassung

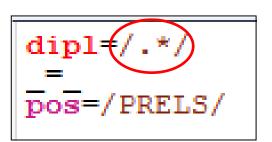


- Suche in ANNIS basiert auf
 - Variablen (Annotationsebenen) und Werten (Kategorien in den Annotationsebenen), z.B.:
 - Suche von exakten Werten, z.B. pos=/ADJA/
 - Suche von Mustern, z.B. pos=/ADJ./
 - Variable-Wert-Paare könne in Relation zu einander gesucht werden, z.B.:
 - Annotationen überlappen sich
 - Annotationen decken einen identischen Bereich ab
 - Variable-Wert-Paare k\u00f6nnen in Verbindung mit Metadaten gesucht werden, z.B.:
 - pos=/ADJA/ & meta::date=/1870/





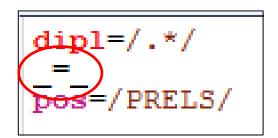
- Ein beliebiges Zeichen
- ? 0 oder 1 Zeichen (des vorherigen Elementes)
- * 0 bis unendlich viele Zeichen (d. vorh. E.)
- + 1 bis unendlich viele Zeichen (d. vorh. E.)
- \\ wörtlich (folgendes Zeichen)
- ! nicht
- [abc] Menge (oder [^abc]=alles außer abc)
- (a|b) a oder b (auch: [ab])
- a{2,3} a 2 bis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



wichtige **Relationen** zwischen VW-Paaren



- . Direkte Präzedenz
- .* Indirekte Präzedenz
- _=_ Identische Abdeckung
- i Inklusion
- _o_ Überlappung
- _l_ linksseitige Überlappung
- _r_ rechtsseitige Überlappung





Vielen Dank!

Anhang



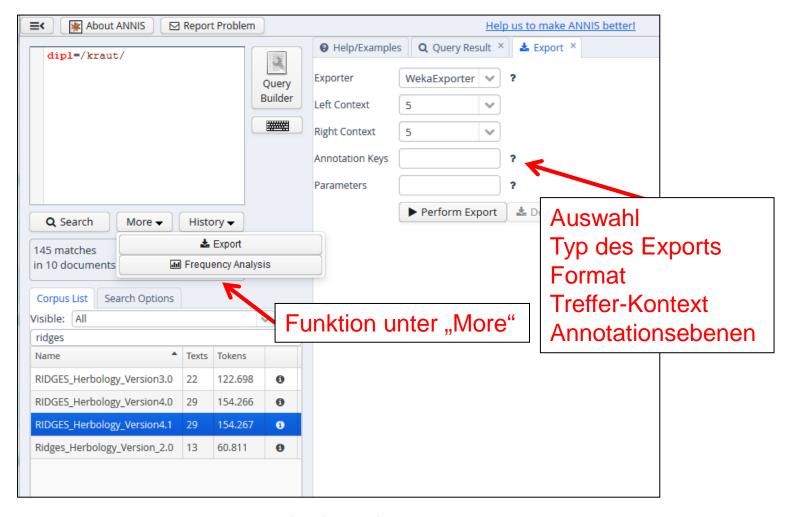
- Interface Export
- Interface Frequency Analysis



Interface Export, Frequenzanalyse

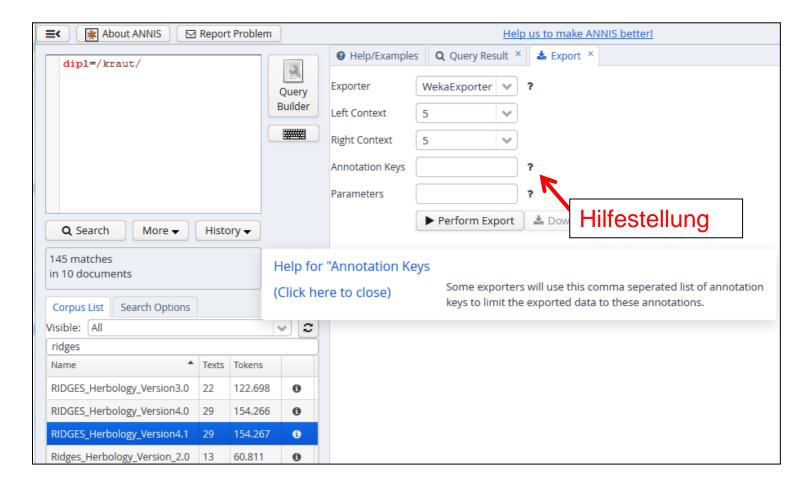












Interface Export von Treffern



Text-Exporter

```
    of the International Brotherhood of [Magicians] Wednesday , October 9
    Magic Month in the United [States] . Wikinews spoke with William
    of the International Brotherhood of [Magicians] , about the current st
    - " Scarne on Card [Tricks] " and " Scarne on
    and " Scarne on Magic [Tricks] " . That started me
```

Grid-Exporter

```
0. tok of the International Brotherhood of Magicians Wednesday
pos IN[1-1] DT[2-2] NP[3-3] NP[4-4] IN[5-5] NPS[6-6] NP[7-7]
cat S[1-6] VP[1-6] NP[1-6] PP[1-6] NP[2-4] PP[5-6] NP[6-6] NP[7-12]
```

CVS-Exporter

```
'11318611','the current state','NP','11318616','current','AJ0','current','JJ'
'11318686','magic','NP','11318688','magic','AJ0','magic','JJ'
'11318757','some basic tricks','NP','11318760','basic','AJ0','basic','JJ'
```

Interface Export von Treffern



Wekaexporter

```
Grelation name
@attribute #1 id string
@attribute #1 span string
@attribute #1 anno const:cat string
@attribute #2 id string
@attribute #2 span string
@attribute #2 anno GUM:claws5 string
@attribute #2 anno GUM:lemma string
@attribute #2 anno GUM:pos string
@data
'11318611', 'the current state', 'NP', '11318616', 'current', 'AJ0', 'current', 'JJ'
'11318686', 'magic', 'NP', '11318688', 'magic', 'AJ0', 'magic', 'JJ'
'11318757', 'some basic tricks', 'NP', '11318760', 'basic', 'AJ0', 'basic', 'JJ'
```





