

ANNIS-Hist: Historische deutschsprachige Korpora in ANNIS

Marcel Bollmann¹, Stefanie Dipper¹, Mario Frank^{1,2}, Julia Krasselt¹, Florian Petran¹, Tom Ruetze²
¹Ruhr-Universität Bochum ²Humboldt-Universität zu Berlin

Einführung

In den letzten Jahren gibt es mehr und mehr Initiativen mit dem Ziel, historische Sprachdaten für die Forschungsgemeinde aufzubereiten und zur Verfügung zu stellen. Dieses Poster stellt eine Reihe von neueren historischen Korpora des Deutschen vor.

Die Korpora werden alle über das Suchtool ANNIS zur Verfügung gestellt. Das ermöglicht eine Suche über alle Korpora hinweg, sodass diachrone Entwicklungen untersucht werden können. Voraussetzung dafür ist die Annotation aller Korpora mit vergleichbaren Tagsets.

HiTS (Historisches TagSet)

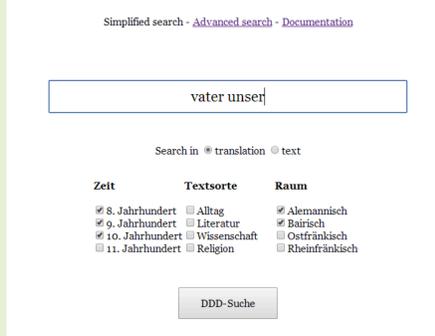
HiTS ist ein gemeinsames Tagset für die Referenzkorpora Altdeutsch, Mittel- und Frühneuhochdeutsch.

- einheitliches Tagset erlaubt diachrone Recherchen
- angelehnt an das STTS
- ergänzende Tags für historische Phänomene, z.B. ADJN für „nachgestelltes attributives Adjektiv“
- zusätzlich: Richtlinien für Tokenisierung (Wortgrenzen)
- zusätzlich: Unterscheidung zwischen Wortart des Lemmas und des Belegs

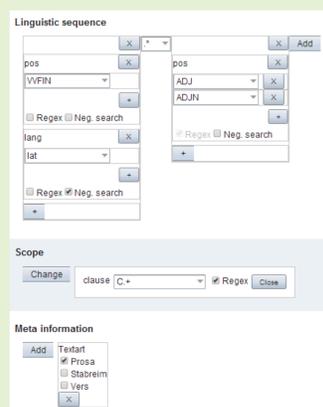
Beispiel:

- inti sie **ni** / **PTK** > **PTKNEG** quedent imo **niouuiht** / **PI** > **PNEG** „und sie nicht sagen ihm Nichts“ = „sie sagen ihm nichts“ (Ahd.)
- swer **niht** / **PI** > **PTKNEG** gloubet, der ist iu verteiltet „wer nicht glaubt, der ist schon verurteilt“ (Späthd.)

Deutsch Diachron Digital



Vereinfachtes Such-Interface von ANNIS



Beispiel-Anfrage im QueryBuilder...

ANNIS

ANNIS ist ein Korpus-Suchtool, das an der HU Berlin entwickelt wird (Zeldes et al., 2009). Das Tool eignet sich insbesondere für Daten, die auf verschiedenen Ebenen mit verschiedenen Typen von Annotationen annotiert sind. Damit unterstützt ANNIS die Idee, dass Korpora sukzessive erweitert und mit weiteren Annotationen angereichert werden können.

Merkmale von ANNIS

1. Visualisierung
 - verschiedene Darstellungsmöglichkeiten, z.B. Tabellenansicht, Graphen, Diskursansicht
 - einzelne Ebenen können verborgen werden
 - Treffer werden farblich hervorgehoben
2. Suche in den Daten
 - AnnisQL: zugrundeliegende formale Anfragesprache
 - QueryBuilder: graphische Oberfläche
 - Vereinfachtes Such-Interface für die historischen Korpora

```
pos="VFIN" & lang!="lat"
& pos=/(ADJ)|(ADJN)/
& clause = /C.+/
& #1_=#2
& #2 .* #3
& #4_i_#1
& #4_i_#3
& meta::Textart = "Prosa"
```

- Anfragen werden in AnnisQL übersetzt und können dort weiter verfeinert werden

Referenzkorpus Altdeutsch (750–1050)

Donhauser (HU Berlin), Gippert (Frankfurt), Lühr (Jena)

Das Referenzkorpus Altdeutsch erfasst die gesamte althochdeutsche und altniederdeutsche Textüberlieferung mit 650 000 Belegen.

Annotationsebenen:

- Wortarten-, Morphologie- und Lemma-Annotation
- Unterscheidung von Lemma- und Beleg-Wortarten (gemäß HiTS)
- Flache Annotation von Satztypen

Beispielphänomen: Negationspartikel *nī*

- typischerweise klitisch (s. Ebene „edition“)
- eigenes Token auf Annotationsebenen (s. Ebene „ling“)

Annotationsansicht (links), Editionsansicht (oben rechts) und Metadaten (unten rechts) in ANNIS																																																			
Path: DDD-Kleinere_Ahd_Denkmaeler > AB_AltbairischeBeichte																																																			
1 1 / 71																																																			
Truhtin, dir uuirdu ih pigihtik allero minero suntiono enti missatatio, alles des ih io missasprah eddo missateta eddo missadahta, uuorto enti uuercho enti kidancho, des ih kihukku eddo nigahukku des ih uuizzanto kiteta eddo unuizzanto, notak eddo unnotak, slaffanti eddo uuachenti: meinsuartio enti lugino, kiridono enti unrehteru																																																			
<table border="1"> <tr><td>edition</td><td>nigahukku</td></tr> <tr><td>ling</td><td>ni</td><td>gahukku</td></tr> <tr><td>lemma</td><td>nī</td><td>gihuggen</td></tr> <tr><td>translation</td><td>nicht</td><td>(ge)denken, denken an</td></tr> <tr><td>posLemma</td><td>PTK</td><td>VV</td></tr> <tr><td>pos</td><td>PTKNEG</td><td>VFIN</td></tr> <tr><td>inflectionClassLemma</td><td></td><td>WK1B</td></tr> <tr><td>inflectionClass</td><td></td><td>WK1B</td></tr> <tr><td>inflection</td><td></td><td>IND_PRES_SG_1</td></tr> <tr><td>clause</td><td>CF_I_Att</td><td></td></tr> <tr><td>line</td><td>3</td><td>4</td></tr> </table>	edition	nigahukku	ling	ni	gahukku	lemma	nī	gihuggen	translation	nicht	(ge)denken, denken an	posLemma	PTK	VV	pos	PTKNEG	VFIN	inflectionClassLemma		WK1B	inflectionClass		WK1B	inflection		IND_PRES_SG_1	clause	CF_I_Att		line	3	4	<table border="1"> <thead> <tr> <th>Name</th> <th>Value</th> </tr> </thead> <tbody> <tr><td>Entstehungszeit</td><td>9,1</td></tr> <tr><td>Referenz</td><td>http://www.paderborner-reperitorium.de/4115</td></tr> <tr><td>Sprache</td><td>ahd.</td></tr> <tr><td>Sprachgebiet</td><td>obd.</td></tr> <tr><td>Sprachlandschaft</td><td>bair.</td></tr> <tr><td>Text</td><td>Altbairische Beichte</td></tr> <tr><td>Textart</td><td>Prosa</td></tr> <tr><td>Textbereich</td><td>Religion</td></tr> </tbody> </table>	Name	Value	Entstehungszeit	9,1	Referenz	http://www.paderborner-reperitorium.de/4115	Sprache	ahd.	Sprachgebiet	obd.	Sprachlandschaft	bair.	Text	Altbairische Beichte	Textart	Prosa	Textbereich	Religion
edition	nigahukku																																																		
ling	ni	gahukku																																																	
lemma	nī	gihuggen																																																	
translation	nicht	(ge)denken, denken an																																																	
posLemma	PTK	VV																																																	
pos	PTKNEG	VFIN																																																	
inflectionClassLemma		WK1B																																																	
inflectionClass		WK1B																																																	
inflection		IND_PRES_SG_1																																																	
clause	CF_I_Att																																																		
line	3	4																																																	
Name	Value																																																		
Entstehungszeit	9,1																																																		
Referenz	http://www.paderborner-reperitorium.de/4115																																																		
Sprache	ahd.																																																		
Sprachgebiet	obd.																																																		
Sprachlandschaft	bair.																																																		
Text	Altbairische Beichte																																																		
Textart	Prosa																																																		
Textbereich	Religion																																																		

Annotationsansicht (links), Editionsansicht (oben rechts) und Metadaten (unten rechts) in ANNIS

rate	fi	bat	daz	fi	liezen	uallin	den	unrechten	willē	uon	ir
rāt	ēr	biten	dazz	ēr	lāzen	vallen	dēr	un-rēht	wille	von	ir(e)
1007	1008						1009				1010
rate	her	bat	daz	fi	liezen	uallin	den	unrechten	willen	uon	ir
Akk_PI	Masc_Nom_Sg_3	Ind_Past_Sg_3	--	*_Nom_PI_3	*_Past_PL_3		Masc_Akk_Sg	Pos_Masc_Akk_Sg_*	Akk_Sg		Neut_Dat_Sg_0
0											
p1											
NA	PPER	VFIN	KOUS	PPER	VFIN	VVINF	DDART	ADJA	NA	APPR	DPOSA
[ra]je	h'	bat	daz	<fi>	[liezen]	[uallin]	[den]	[unrechten]	[w]ille	uon	ir

Beispielergbnis aus dem Mittelhochdeutsch-Korpus bei Suche nach VFIN gefolgt von VVINF (Suchanfrage in AnnisQL: `pos="VFIN" & pos="VVINF" & #1 . #2`)

Referenzkorpus Mittelhochdeutsch (1050–1350)

Wegera, Dipper (Bochum); Wich-Reif, Klein (Bonn)

Das Referenzkorpus Mittelhochdeutsch enthält eine strukturierte Auswahl der Textüberlieferung mit etwa 2,1 Mio. annotierten Wortformen.

Annotationsebenen:

- Wortarten-, Morphologie- und Lemma-Annotation
- Unterscheidung von Lemma- und Beleg-Wortarten (nach HiTS)

Beispielphänomen: Abfolge VFIN/VVINF im Verbalkomplex

- z.B. *dass sie fallen ließen* vs. *dass sie ließen fallen*
- AnnisQL erlaubt Suche nach konkreter Abfolge von Wortarten-Tags

Anselm-Korpus (Frühneuhochdeutsch, 14.–16. Jh.)

Dipper, Schultz-Balluff, Wegera (Bochum)

Das Anselm-Korpus umfasst 24 frühneuhochdeutsche Überlieferungsträger des Passionstraktats „St. Anselmi Fragen an Maria“. Die Überlieferungen variieren in Länge, Textsorte (Vers vs. Prosa) und dialektalen Merkmalen.

Annotationsebenen:

- Wortarten-Annotation nach STTS
- Normalisierung: streng wortnahe Anpassung der Form an moderne Schreibung
- Modernisierung: Anpassung der Formen an moderne Grammatik und Semantik:
 1. Anpassung an moderne Flektion („fлект“)
 2. Lexikalische Anpassung extinkter Formen („lex“)
 3. Anpassung bzgl. semantischem Wandel der historischen Formen („sem“)

Ausschnitt aus dem Anselm-Korpus in ANNIS; historische Wortformen können in UTF-Darstellung oder simplifizierter Darstellung (ersetzt z.B. f durch s) gesucht werden	
Path: Me1_Meik > Me1_Meik	
alle creatur	die got pēchueff douon So mag ich noch en fchol nymer wainenn vnd douon fo will ich von anengeng
PIAT NN	PRELS NN VFIN PAV ADV VMFIN PPER ADV PTKNEG VMFIN ADV VVINF KON PAV ADV VMFIN PPER APPR NN
flekt	sem lex
alle creatur	die got peschueff douon So mag ich noch en schol nymer wainenn vnd douon so will ich von anengeng
alle Kreatur	die Gott beschuf davon so mag ich noch en soll nimmer weinen und davon so will ich vom Aneganc
alle Kreaturen	die Gott schuf davon so kann ich noch soll nimmer weinen und davon so will ich vom Anfang
inline (grid)	
paula	
paula text	

Ausschnitt aus dem Anselm-Korpus in ANNIS; historische Wortformen können in UTF-Darstellung oder simplifizierter Darstellung (ersetzt z.B. f durch s) gesucht werden

Beispielphänomen: Regelmäßige graphematische Ersetzungen

- z.B. *v* → *u* (wie in *vnd* → *und*) oder *u* → *v* (wie in *douon* → *davon*)
- Normalisierungsebene erlaubt Vergleich historischer und moderner Schreibvarianten