

TEITOK



**COMBINING LANGUAGE AND LINGUISTIC
INFORMATION WITHOUT COMPROMISE**

Language vs. Linguistics



- **Linguistic information (ANNIS)**
 - Tokenization, POS, lemma, dependencies, syntax, semantics
 - Normalizations, sentences
- **Language information (TEI)**
 - Lines, pages, paragraphs
 - Bold, italics, colour, typesettings
 - Hands, additions, deletions
- **Document information (Dublin Core)**
 - Facsimile images
 - Year, library, author, etc.

Traditional Method



- Annotated Corpora are text based
- Step 1 is to clean the text
 - Throw away any non-linguistic information
- Step 2 is to tokenize
 - Verticalizing text one word per line
 - Removes information about inter-word spacing
 - Often also splitting contractions

Facsimile Image

viij genueyn.

Don dem S.

Wer eyne hübsche farbe wil haben.

Btonica. Trinck wein ab Betonien/so wird dir ein gute farb spricht Plinius. Wer sie bey ihm trage / dem mag keyne zauberey schaden. Es ist auch gut fur gifft/vnd wer einen bösen Magen hat/Leber vnd Milz/der mag trincken ab dem kraut / also / das darunder gemischt werd ein wenig essig vnd honig/diss also getruncken es hilfft/es ist auch gut denen die blut speyen.

Wer ein gut gedechnis wil haben.

Buglossa. Wer Ochsen zungen kraut beysset yn wein/ vnd darnach trinckt/der gewint ein gut gedechnis. Es sterckt auch das hertz/vnd macht gut blut/vn heilet auch das hertz gesper. Den safft getruncken mit warmen wass

Example from the RIDGES corpus

Verticalized text



A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	
dipI	clean	norm	pos	lemma	pos_klein	Verpositi	KOUS_Ser	Nebensatz	komp	komp_orth	prot	attr_gen	strD	erlaeuteru	unclear	atLeast	atMost	interpretat	hi_rend	typeface	note
Von	Von	Von	APPR	von	APPR															gothic	
dem	dem	dem	ART	die	ART															gothic	
B.	B.	B.	NN	B.	N															gothic	
Wer	Wer	Wer	PWS	wer	N															gothic	
eyne	eyne	eine	ART	eine	ART															gothic	
hubiche	hubische	hubische	ADJA	hubisch	ADJ															gothic	
farbe	farbe	Farbe	NN	Farbe	N															gothic	
wil	wil	will	VMFIN	wollen	VFIN															gothic	
haben	haben	haben	VAINF	haben	VINF															gothic	
.	.	.	\$.	.	ZEICHEN															gothic	
BEtonica	BEtonica	Betonica	NE	<unknown>	N															iniCap	gothic
.	.	.	\$.	.	ZEICHEN															gothic	
Trinck	Trinck	Trink	VVIMP	trinken	VIMP															gothic	
wein	wein	Wein	NN	Wein	N															gothic	
ab	ab	ab	APPR	ab	APPR															gothic	
Betonien	Betonien	Betonien	NN	<unknown>	N															gothic	
/	/	/	\$({	/	ZEICHEN															gothic	
fo	so	so	ADV	so	ADV															gothic	
wird	wird	wird	VAFIN	werden	VFIN															gothic	
dir	dir	dir	PPER	du	N															gothic	
ein	ein	eine	ART	eine	ART															gothic	
gute	gute	gute	ADJA	gut	ADJ															gothic	
farb	farb	Farbe	NN	Farbe	N															gothic	
þricht	spricht	spricht	VVFIN	sprechen	VFIN															gothic	
Plinius	Plinius	Plinius	NE	<unknown>	N															gothic	
.	.	.	\$.	.	ZEICHEN															gothic	
Wer	Wer	Wer	PWS	wer	N															gothic	
fie	sie	sie	PPER	sie	N															gothic	
bey	bey	bei	APPR	bei	APPR															gothic	
ihm	ihm	ihm	PPER	er	N															gothic	
trage	trage	trage	VVFIN	tragen	VFIN															gothic	
/	/	/	\$({	/	ZEICHEN															gothic	
dem	dem	dem	PRELS	die	N															gothic	
mag	mag	mag	VMFIN	mögen	VFIN															gothic	
keyne	keyne	keine	PIAT	keine	ART															gothic	
zauberey	zauberey	Zauberei	NN	Zauberei	N															gothic	

Final version (Excel) of the verticalized text with columns

ANNIS - Text



Help us to make ANNIS better!

not logged in [Login](#)

dialekt="alemannisch"

[Query Builder](#)

[Keyboard](#)

[Search](#) [More ▾](#) [History ▾](#)

Valid query, click on "Search" to start searching.

[Corpus List](#) [Search Options](#)

Visible: All

Filter

Text	Tokens	Info	View
GES_Herbology_Version3.0	22	122,698	Info
GES_Herbology_Version4.0	29	154,266	Info
GES_Herbology_Version4.1	29	154,267	Info View
GES_Herbology_Version5.0	36	183,724	Info View
ges_Herbology_Version_2.0	13	60,811	Info
JLTRON_Banana	2	3,782	Info View

RIDGES_Herbology_Version5.0 > ArtzneyBuchleinDerKreutter_1532_Tallat - Visualizer: normalized transcription

Von dem B.

Wer eine hübsche Farbe will haben .

Betonica . Trink Wein ab Betonien / so wird
dir eine gute Farbe spricht Plinius . Wer sie bei
ihm trage / dem mag keine Zauberei schaden .
Es ist auch gut für Gift / und wer einen bösen
Magen hat / Leber und Milz / der mag trinken
ab dem Kraut / also / dass darunter gemischt wird ein
wenig Essig und Honig / dies also getrunken es hilft / es ist
auch gut denen die Blut speien .

Wer ein gutes Gedächtnis will haben .

Buglossa . Wer Ochsenzungenkraut beißt in Wein /
und danach trinkt / der gewinnt ein gutes Gedächtnis . Es
stärkt auch das Herz / und macht gutes Blut / und heilt auch
das Herzgesperr . Den Saft getrunken mit warmem Wasser
hilft für das Geschwelle der Füße .

Was den Durst bemehm .

Berberis . Brauche Pfirsich so benimmt er den Durst /
und stärkt den Magen und die Leber . Item damit geschmiert
den Bauch der Frauen / treibt aus das tote Kind /

auch macht es schwitzen . Pfirsich mit Wasser das Saft
ausgedrückt / und davon morgens genossen / ist gut wider
das Hauptweh / spricht Platearius .

Für Schwindel in dem Haupt .

Borrage . Das Saft von dem Kraut Borrich / misch mit
Zucker und trink es / das hilft für den Schwindel im Haupt .

Linguistics Only



- Corpora mostly cannot be rendered as “text”
 - Example before still has linebreaks, making it better to read
 - Full-text representation often not even offered
- Good for linguistics
 - Statistical analysis
 - Complex queries
 - Grammatical analysis
- Bad for most other things
 - Not friendly (or even possible) to read
 - No link to facsimile images

TEITOK



- TEITOK starts from TEI/XML
 - Each file is in full TEI/XML format
- Annotation is added inline
 - Tokens are added non-destructively to the XML
 - Annotations are attributes over nodes
- Extralinguistic rendering
 - TEI encodings (such as dropcaps) rendered as CSS
 - Images rendered next to document (from <pb/>)
- Graphical display in browser
 - Examples to follow

RIDGES

[index](#)

Folio 17

>

Von dem B.

Wer eyne hübsche farbe wil haben.

Betonica. Trinck wein ab Betonien / so wird dir ein gute farb spricht Plinius. Wer sie bey ihm trage / dem mag keyne zauberey schaden.

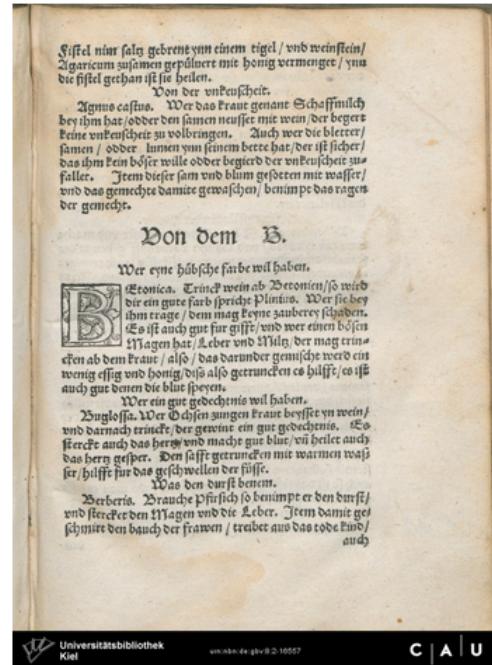
Es ist auch gut fur gifft / vnd wer einen boffen Magen hat / Leber vnd Miltz / der mag trinken ab dem kraut / also / das darunder gemischt werd ein wenig effig vnd honig / diſs also getruncken es hilfft / es ist auch gut denen die blut speyen.

Wer ein gut gedechnis wil haben.

Buglossa. Wer Ochsen zungen kraut beyſſet yn wein / vnd darnach trinckt / der gewint ein gut gedechnis. Es ſterckt auch das hertz / vnd macht gut blut / v heilet auch das hertz gesper. Den ſafft getruncken mit warmen waſſer / hilfft fur das geschwellen der fuſſe.

Was den durft benem.

Berberis. Brauche Pfirsich fo benimpt er den durft / vnd ſtercket den Magen vnd die Leber. Item damit geſchmiert den bauch der frawen / treibet aus das tote kind /



Universitätsbibliothek
Kiel

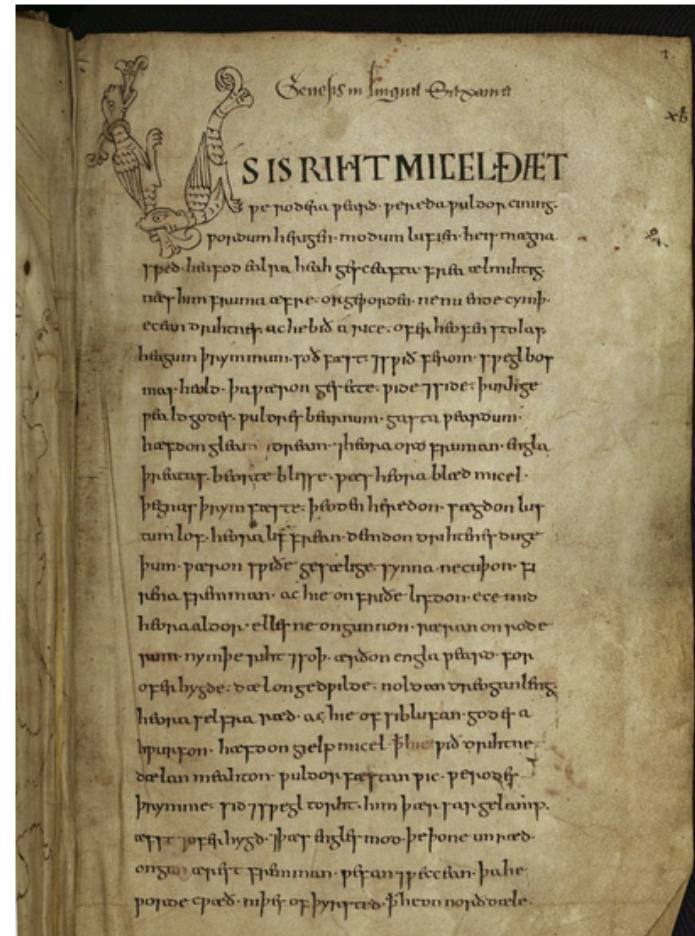
C | A | U

Junius 11



US IS RIHT MICEL ÐÆT

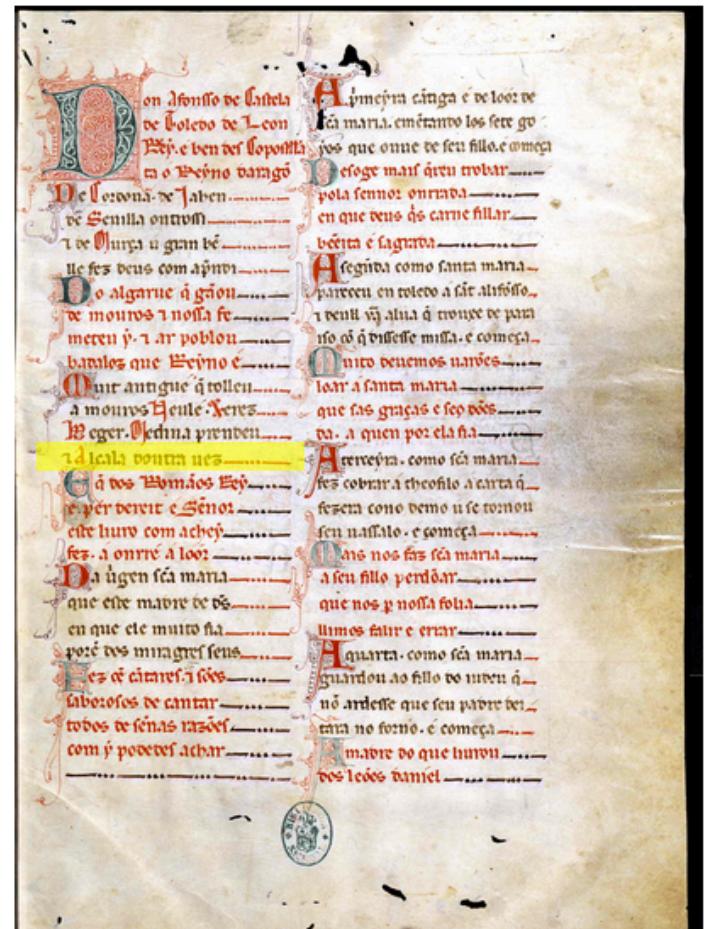
pe rodera peard · pereda puldorcining ·
þordum herigen · modum lufien · he iſ mægna
ſped · heafod ealra heah gerceafta · frea ælmihtig ·
nær him fruma æfre · or geporden · ne nu ende cymb þ ·
ecean drihtner · ac he bið a rice · ofer heofen ƿtolar ·
heagum þrymmum · ƿoðfært 7 ƿið ferom ƿeglbor
mar heold · þa pæron gerette · pide 7 ride · þurh ge
peald goder · puldrer bearnum · gaſta peardum ·
hæfdon gleam 7 dream · 7 heora ord fruman · engla
þreatar · beorhte bliſſe · pær heora blæd micel ·
þegnar þrymfæſte · þeoden heredon · rægdon lur
tum lof · heora liffrean · demdon drihtener dugē
bum · pæron ƿiðe gerælige · ƿynna ne cuþon · fi
rena fremman · ac hie on friðe lifdon · ece mid
heora aldon · eller ne ongunnon · ræran on rode
rum · nymþe riht 7 ƿob · ærdon engla þlaflo · for
oferhygde · dæl on gedpilde · noldan dreogan leng ·
heora ƿelfra ræd · ac hie of ƿiblufan · goder a
hpurfon · hæfdon gielp micel · þa hie pið drihtne ·



Cantigas de Santa Maria



Don Afonso de Castela
de Toledo de Leon
Rey e ben del Copostela
ta o Reyno daragō
De Cordoua • de Iahlen
de Seuilla outrossi
γ de Murça u gran bē
lle fez deus com aþndi
Do algarue q gāou
de mouros γ noſſa fe
meteu y • γ ar poblou
badaloz que Reyno e
Muit antigue q tolleu
a mouros Neule • Xerez
Beger • Medina prendeu
e **A**lcala doutra uez
E q dos Romãos Rey
e per dereit e Señor
este liuro com achey
fez a onrre a loor
Nra virgen fia maria



Zographensis



ԵՍՏՐԱՔՅԱ ՅԾՆԻ

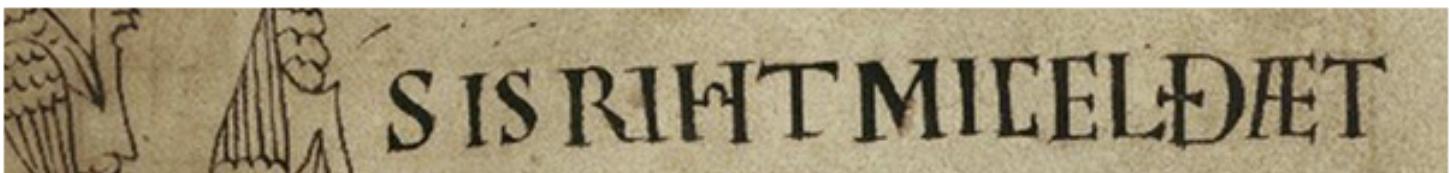
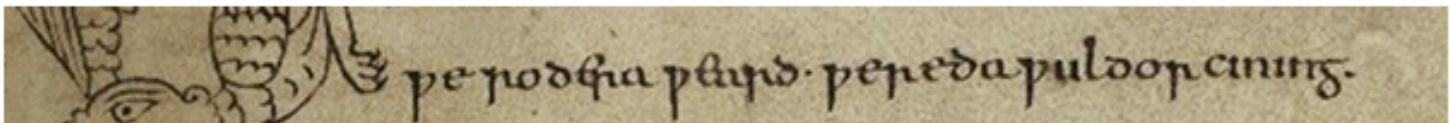
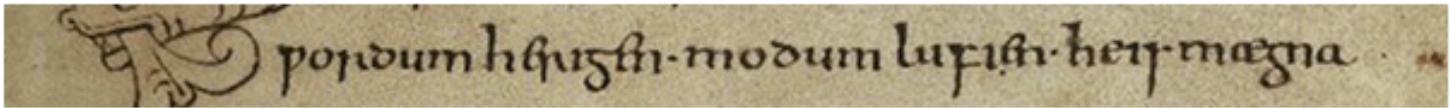
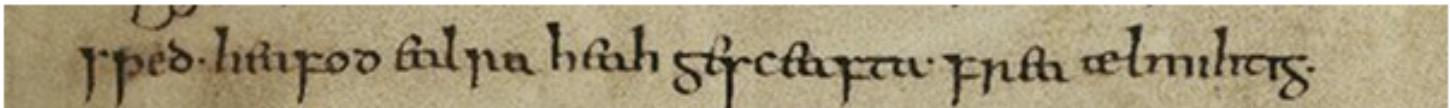
ՀՏԵՔՆԻ :

ԽԱՇՎԵԱ ԵՍՏՐԱՔՅԱ . ԴՉ ԵՄԴ
ԶԲԴ ԸԱՋԱ . ՃԱԵԿԵ Յ
ԶԹԱ ՐԱԴԻՔՅ ՄՅ ՐԵԵԵՎԱԽԵՅ .
ԶԵ ՇԽԱ ՐԵԶԵՔՅՈՒՆ ՆԱՔՆ ՀՅԵՅ .
ՐԵԱԽԵ ԽԵՎԵՅՅ ՀԵՎՅՅ . ԴԿԵ
ՊԵԿԵԵՄՅՅ ՐԵԿԵՄՅ ԾՄԵԴ . ՐԵԱ
ԽԵ ԾԵԱԿՅ . ԿԵԴԻՔՅ ՄԵՐՄՅ
ՄԹԵԿՅ ՄՅ ՐԵՎՅՄՅՔՅ . ՊԵԿ
ԾԵՄԵԴՅ ՐԵԿԵՄՅ ԽԵՔՅ . ՐԵԱՄ
ԾՄԵՎՅՅ ԶԹԱԲԵԿ ՅԿԵ . ԾԵԿ
ԶԹԱ ԴԵԴԻՔՅ ԽԵՎՅՄՅ ՄՅ ՐԵՎՅՄՅ
ՔՅ . Դ ՐԵԵՐՄՅՄՆԵՅ ԽԵՎՅՄՅ
ՔՅՅ ՐԵԱՆԻԴՔՅ . ՄՅ ՅԾՆԵՐՄՅ
ՄԹԵՔՅՅ ԽԵՎՅՅ . Դ ՔՅՅ



Junius 11 – Manuscript lines

Folio 1

[1]	 US IS RIHT MICEL ÐÆT
[2]	 pe rodera peard · pereda puldorcing ·
[3]	 pordum herigen · modum lufien · he iſ mægna
[4]	 rned · heafod ealra heah cerceafta · frea ælmihtig ·

- Transcription can be aligned to manuscript line
- Good for verifying transcriptions against the original

Inline Tokenization



- TEI can be complex
 - 3 options to add information
- Export only words
 - Creates unrelated documents
- Stand-off annotation –
 - Typically position-based in the text: chars. 12-16 “Buch”
 - Difficult with XML, refutes robustness, and not editable
- Inline annotation
 - Create XML nodes for tokens
 - “Das Buch hier” => <token>Buch</token>
 - Unmodified TEI

XML Attributes



- Annotations marked over nodes
 - <tok pos="Noun">Buch</tok>
 - TEI: <w ana="N">Buch</w>
- White-space sensitive
 - Spaces in the text are meaningful
 - <c> </c>
- Inline verticalization
 - When looking only at <tok/> we get a list of tokens
 - Easy to convert to/from .vrt format

Multiple Orthographies



- Normalization typically required in historic corpus
 - Original text also needed
- Multiple orthographic layers
 - RIDGES has those as well
- No limit to the amount of forms
 - Original and normalized
 - Diplomatic, modernized, etc.
 - Romanized, transliterated
- Hierarchical
 - Dependent layers, only filled when needed

Graphical User Interface



- XML of TEITOK not revolutionary
 - Very complex
 - <pb/> Von dem B.
 - <pb facs="p17.jpg" n="17" id="e-1"/> <p> <lb id="e-2"/> <tok form="Von" nform="Von" pos="APPR" lemma="von" id="w-1">Von</tok> <tok form="dem" nform="dem" pos="ART" lemma="die" id="w-2">dem</tok> <tok form="B." nform="B." pos="NN" lemma="B." id="w-3">B.</tok> <lb id="e-3"/> </p>
- Built to allow using complex XML
 - Simple HTML forms wherever possible

CQP Corpus



- A TEITOK corpus consists of XML files
 - Not easy to search
- Automatic exporting to CQP
 - Dedicated export module tt-cwb-export
 - Fast enough even for “large” corpora (10 minutes)
- Corpus search as XML index
 - Added a level to CQP for byte-offset in XML file
 - CQL => byte-offset => XML
 - Full XML nodes shown in KWIC

Variable Tokenization



- Tokens are not identical at every level
 - hertz gesper => Hertzgesperr
 - Non-aligned tiers in ANNIS
- Spaces in the appropriate level
 - <tok nform="Hertzgesperr">hertz gesper</tok>
 - <tok nform="Auf dem">Auffm</tok>
- Grammatical words as subtokens
 - <tok>Auffm<dtok form="auf"/><dtok form="dem"/></tok>
 - Reversely usable for MWE <mtok><tok/><tok/></mtok>
- Requires a choice when exporting to CWB

Additional Modules



- Dependency Trees
 - [http://alfclul.clul.ul.pt/teitok/test/index.php?
action=deptree&cid=catconll2009.xml&sid=s-3](http://alfclul.clul.ul.pt/teitok/test/index.php?action=deptree&cid=catconll2009.xml&sid=s-3)
- Syntactic Trees
 - [http://ps.clul.ul.pt/index.php?
action=psdx&cid=CARDSo0001&treeid=tree-1&node=node-8](http://ps.clul.ul.pt/index.php?action=psdx&cid=CARDSo0001&treeid=tree-1&node=node-8)
- Stand-off Annotations (error annotation)
 - [http://alfclul.clul.ul.pt/teitok/learnercorpus/index.php?
action=annotation&annotation=errors&cid=nlo01CVETD.xml
&query=](http://alfclul.clul.ul.pt/teitok/learnercorpus/index.php?action=annotation&annotation=errors&cid=nlo01CVETD.xml&query=)

Additional Modules



- Interlinear Glossed Text
 - [http://90.171.34.41/teitok/typecraft/index.php?
action=igt&id=ada/2175.xml](http://90.171.34.41/teitok/typecraft/index.php?action=igt&id=ada/2175.xml)