



owards machine readable linguistic resources for Medieval Latin

from non-structured editions
to automatic parsing

"Towards machine readable linguistic resources for Medieval Latin – from non-structured editions to automatic parsing"

Hanna-Mari Kupari

DH pizza lunch seminar 14th April, 2023

Aalto University



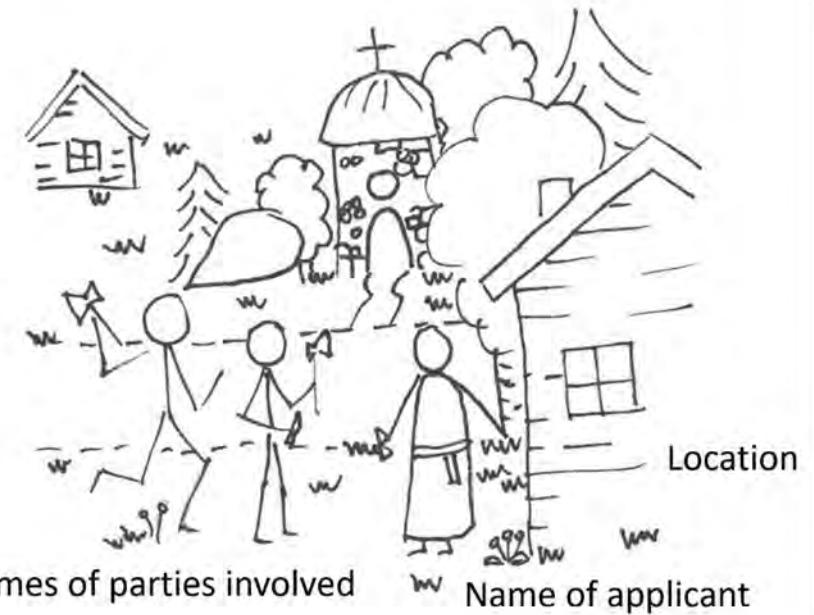
EMIL AALTOSEN SÄÄTIÖ

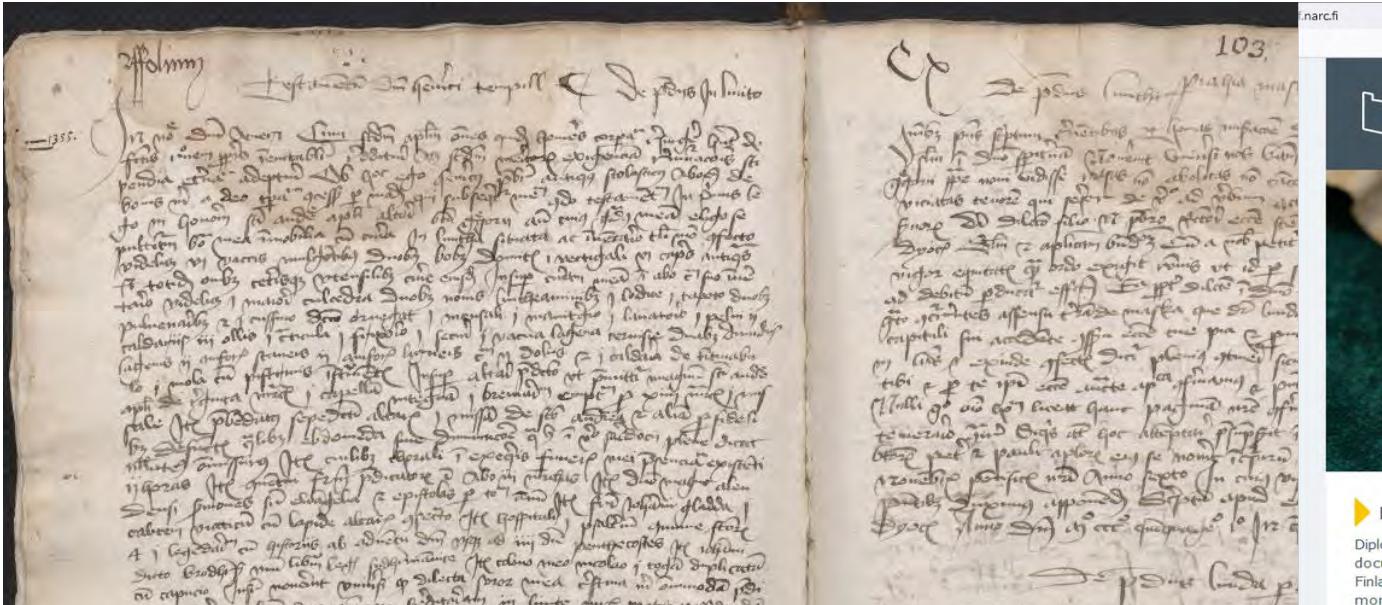
Outline of the presentation

Short background

1. The data
2. The structuring
3. The automatic morpho-syntactic parsing

Discussion





The screenshot shows a Visual Studio Code interface with two open files:

- landscape-analysis.py**: A Python script using ElementTree to parse `dataset_landscape.xml`. The script counts gender and status across documents and creates a pandas DataFrame.
- dataset_landscape.xml**: An XML document containing historical records of individuals (issuers) from the Borgaren database.

Below the code editor, the terminal window displays:

- A warning about NumPy being compiled with Python 3.9.
- A note to check the Python and NumPy versions.
- A reminder to study the documentation for troubleshooting.
- The original error message: `DLL load failed while importing _multiarray_umath: Määritettyä osaa ei löydy.`
- The command `PS C:\Users\Kanna-Mari\Documents\Tohtorius\Kielimaisemata Lämberg\landscapes> []`.

At the bottom, the status bar shows: In 21, Col 1 | Spaces: 2 | UTF-8 | CR LF | Python 3.9.13 (base) | 161 lines.

The screenshot shows the homepage of the Diplomatarium Fennicum. At the top left is the URL 'narc.fi'. The top right features a search bar with a magnifying glass icon and language links 'FI SV EN'. Below the header is a dark blue banner with the text 'DIPLOMATARIUM FENNICUM' and a stylized book icon. To the right of the banner are three buttons: 'SEARCH' with a magnifying glass, 'HELP' with a question mark, and 'MENU' with a list icon. The main content area displays a historical document on parchment with handwritten Latin text. A large, textured seal or stamp is visible on the left side of the document. At the bottom right of the page is a dark box containing the text 'DIPLOMATARIUM FENNICUM'.



Short bio

- High school Latin - 2003
 - BA: Latin 2014 – late antiquity, quantitative methods and gendered language use
 - MA: Latin 2019 – Medieval Latin, basic digital methods and the penitentiary documents: expressions of violence and death
 - PhD studies: digital linguistics major – university of Turku

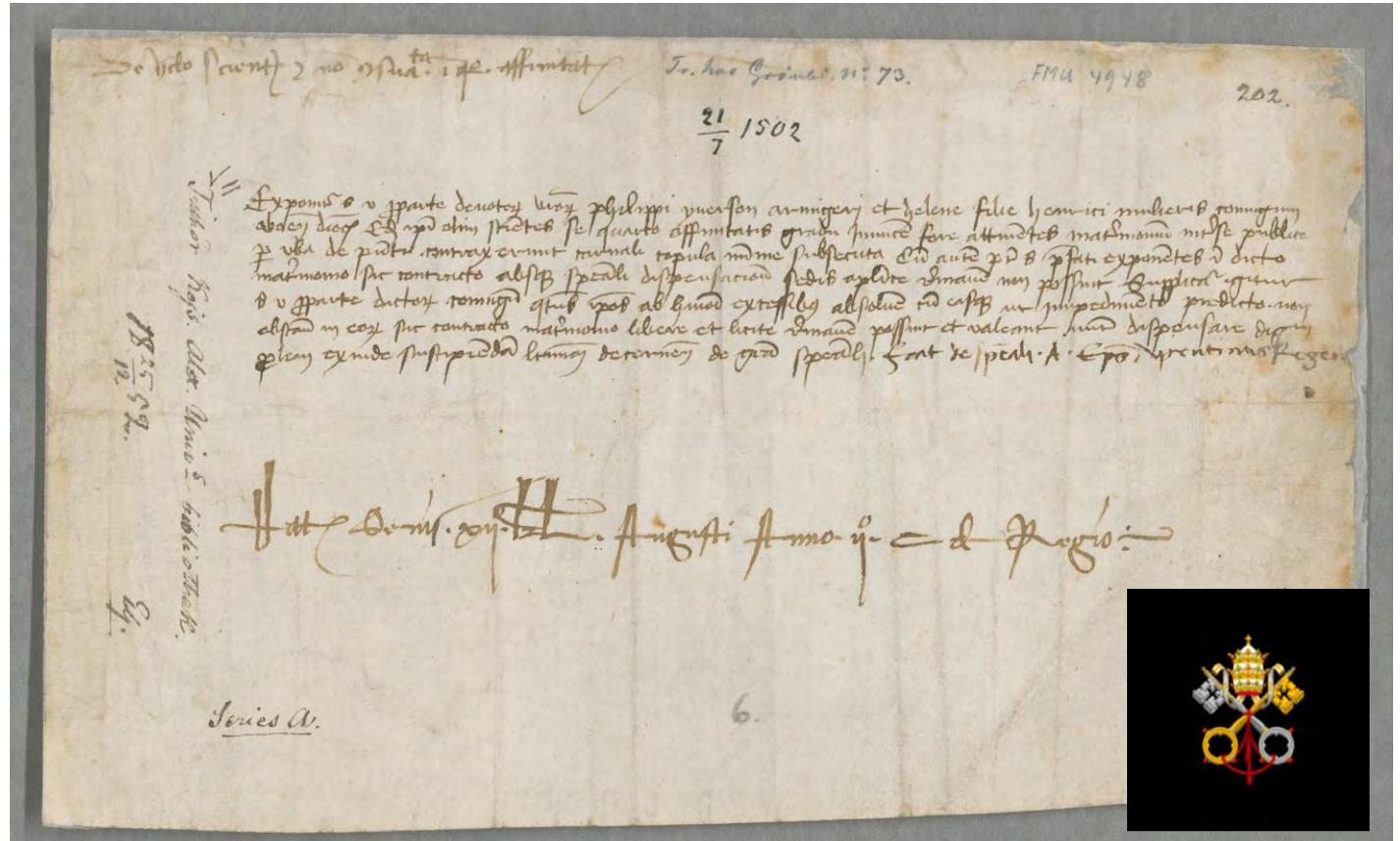


Section 1

The data

Motivation

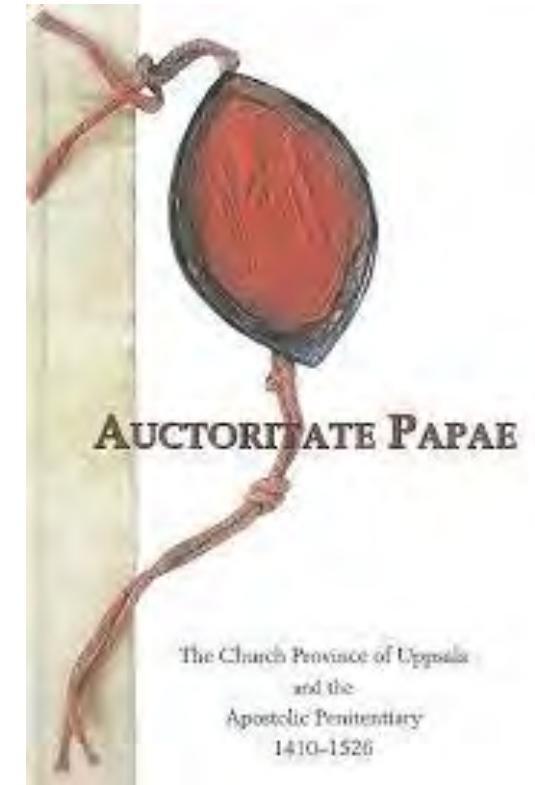
- Professor Kirsi Salonen
- Local history
- Historical linguistics
- Insightful and fun
- Novel research



<http://df.narc.fi/document/4948>

Penitentiary documents

- From the Apostolic See
- Copybooks survive and edited according to modern principles for historians using quantitative methods
- The Finnish documents ->



Exponitur sanctitati vestre pro parte devotorum vestrorum Philippi Yverson armigeri et Helene filie//
Henrici mulieris coniugum Aboensis diocesis, quod ipsi olim scientes se quarto affinitatis gradu invicem fore
attinentes matrimonium inter se publice per verba de presenti contraxerunt carnali copula minime
subsecuta. /

Exponitur sanctitati vestre pro parte devotorum vestrorum Philippi Yverson armigeri et Helene filie//
Henrici mulieris coniugum Aboensis diocesis, quod ipsi olim scientes se quarto affinitatis gradu invicem fore
attinentes matrimonium inter se publice per verba de presenti contraxerunt carnali copula minime
subsecuta. /

The context of language use with the penitentiary documents



The situation – e.g. violent assault



description



+
Linguistic
features

Section 2

The structuring

From Data to Capta



<https://www.shopsy.in>

Johanna Drucker

- *Humanities Approaches to Graphical Display:*
- Capta is “taken” actively while data is assumed to be a “given” able to be recorded and observed.
- From this distinction, a world of differences arises
- knowledge is constructed, taken, not simply given as a natural representation of pre-existing fact
- <http://www.digitalhumanities.org/dhq/vol/5/1/000091/000091.html>

The structuring

- What works for human consumption is really difficult for machine reading
- The end point of an editor's work is where I start – pdf's
- By the means of regex and manual annotation
- Balancing between accuracy and number of categories



TEI-XML markdown

442

30.9 1521

Rome

Henricus Johannis, a Birgittine brother in the monastery of Naantali (Nådendal) in the diocese of Turku and a priest and Master of Arts, is forty years old and suffers from insomnia, weak eyes, constipation, hemorrhoids and other diseases. Because of this, he cannot endure the duties of the monastery and the order, but wishes to be

441,14 et...18 obstruenda addidi secundum n. 358,14–17, etc cod. | 18 etc] scil. humiliter eidem sanctitatis vestre pro parte ipsius exponentis vel sim.; cfr e.g. n. 119,19; 278,13; 341,8. [24 Abensi] i.e. Aboensis | 442. Anno nono Leonis pape x in marg. sup.; i Kalendas Octobris in marg. sin.; Aboensis in marg. dext.; Dispensatio ut infra; nomen procuratoris Cencis vel sim.; turonenses xi; residua gratis intulit procuratoris ante annot. fol. 534r; Rome apud Sanctum Petrum in marg. sup. fol. 534v.

463

Sara Risberg

exempted from the duties of the choir for as long as he is afflicted with these diseases.
The regent Mercurius grants Henricus his request.

Henricus Johannis monasterii Vallisgratiae ordinis sancti Salvatoris Aboensis diocesis professor exponit, quod ipse, qui presbyter et in artibus magister ac quadragenarius existit, propter insomnepitatem et debilitatem oculorum suorum ventrisque constipationem ac 5 hemorroidarum fluxum et alias sui corporis infirmitates, quas sepius patitur, austерitas et onera regularia monasterii et ordinis predictorum suffere posse non sperat. Et propterea, ut eisdem infirmitatibus consulere possit, cuperet ab eisdem ac oneribus chori eximi, quod sibi

68:534r

```
trankit_custom_model.py   Tex Translations  autoritate_teksti.xml X
Data > autoritate_teksti.xml > autoritate > document > h2
7460 emulorum obstruenda </i>supplicatur etc., quatenus, si iudici infrascripto
7461 canonice constiterit de iusto bello et aliis premissis, ipsum oratorem
7462 nullam propter premissa inhabilitatis vel irregularitatis maculam sive
7463 notam contraxisse sed premissis non obstantibus in suis ordinibus []
7464 etiam in altaris ministerio [] ministrare libere et licite posse nuntiari et
7465 declarari mandare dignemini prout in forma. Fiat in forma. Mercurius
7466 regens. Committatur episcopo Abuensi vel eius vicario in spiritualibus.
7467 Fiat. Mercurius.
7468 </document>
7469
7470 <document>
7471 <h2 num="442" niputus="n" several_witnesses="n">442 <date when="1521-09-30">30.9 1521</date> Rome</h2>
7472 Henricus Johannis monasterii Vallisgratiae ordinis sancti Salvatoris
7473 Aboensis diocesis professor exponit, quod ipse, qui presbyter et in
7474 artibus magister ac quadragenarius existit, propter insomnepitatem et
7475 debilitatem oculorum suorum ventrisque constipationem ac
7476 hemorroidarum fluxum et alias sui corporis infirmitates, quas sepius
7477 patitur, austерitas et onera regularia monasterii et ordinis predictorum
7478 suffere posse non sperat, et propterea, ut eisdem infirmitatibus consulere
7479 possit, cuperet ab eisdem ac oneribus chori eximi, quod sibi
7480 permitti dubitat sede apostolica desuper inconsulta. Supplicatur igitur
7481 humiliter pro parte dicti oratoris, quatenus se ab eisdem austерitatibus
7482 ac intercessione chorii tam in diurnis quam nocturnis horis infirmitatibus
7483 huicmodi durantibus, dummodo propterea divinum officium recitare
7484 non omittat, veris existentibus premissis, cuius super hoc parte licentia
7485 minime requisita eximi et sibi indulgeri secundum desuper misericorditer
7486 dispensari non obstantibus constitutionibus et ordinationibus
7487 apostolicis ac tam provincialibus quam synodalibus necnon monasterii
7488 et ordinis predictorum statutis et consuetudinibus iuramento etc.
7489 roboratis privilegiis quoque indultis et litteris apostolicis monasterio et
7490 ordini huiusmodi sub quibusvis verborum formis et clausulis etiam
7491 derogatoriis derogatoriis, fortioribus, efficacioribus et insolitis contra
7492 premissa forsan concessis et imposterum concedendis, etiam si de illis
7493 eorumque totis tenoribus pro illorum sufficienti derogatione specialis,
7494 specifica, expressa et individua ac de verbo ad verbum, non autem etc.
7495 habenda foret et alias, cuiuscumque tenoris existant, ceterisque contrariais
7496 quibuscumque dignemini de gratia speciali et expresso. Fiat de
7497 speciali et expresso. Mercurius regens.
7498 </document>
7499
7500 <document>
7501 <h2 num="443" niputus="n" several_witnesses="n">443 <date when="1522-06-28">28.6 1522</date> Rome</h2>
7502 Andreas Gladh diaconus Arosiensis diocesis in artibus magister exponit,
7503 quod ipse ex magno devotionis fervore cupit ad presbyteratus ordinem
7504 promoveri et promotus in illo ac aliis iam per eum susceptis ordinibus []
7505 etiam in altaris ministerio [] ministrare ac beneficia ecclesiastica cum
7506 cura ac sine cura invicem tamen compatibilia sibi in futurum canonicę
7507 conferenda recipere et retinere posse. Sed quia quandam maculam in
7508 oculo suo dextro, ex quo forsan nichil videt, patitur, desiderium suum
```

Many opportunities of xml



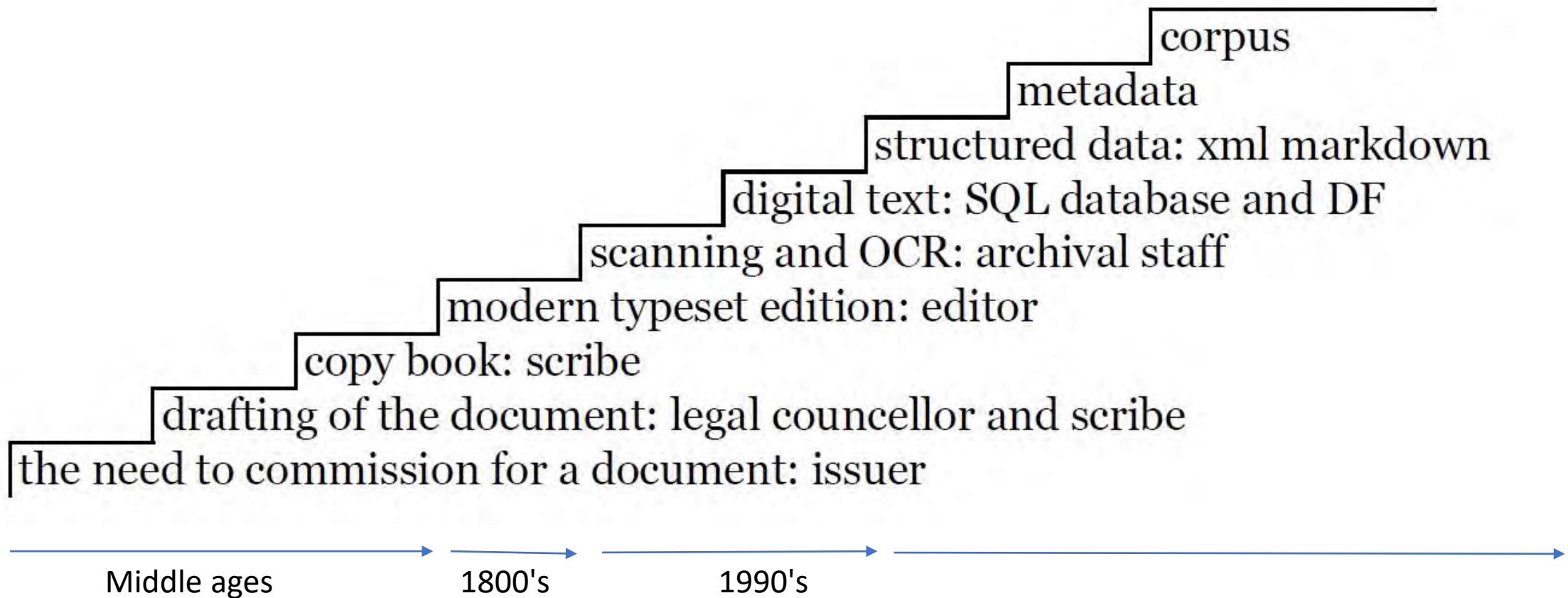
<https://yle.fi/a/3-12159318>

Linguistics landscapes: the use of landmarks in Turku copybooks

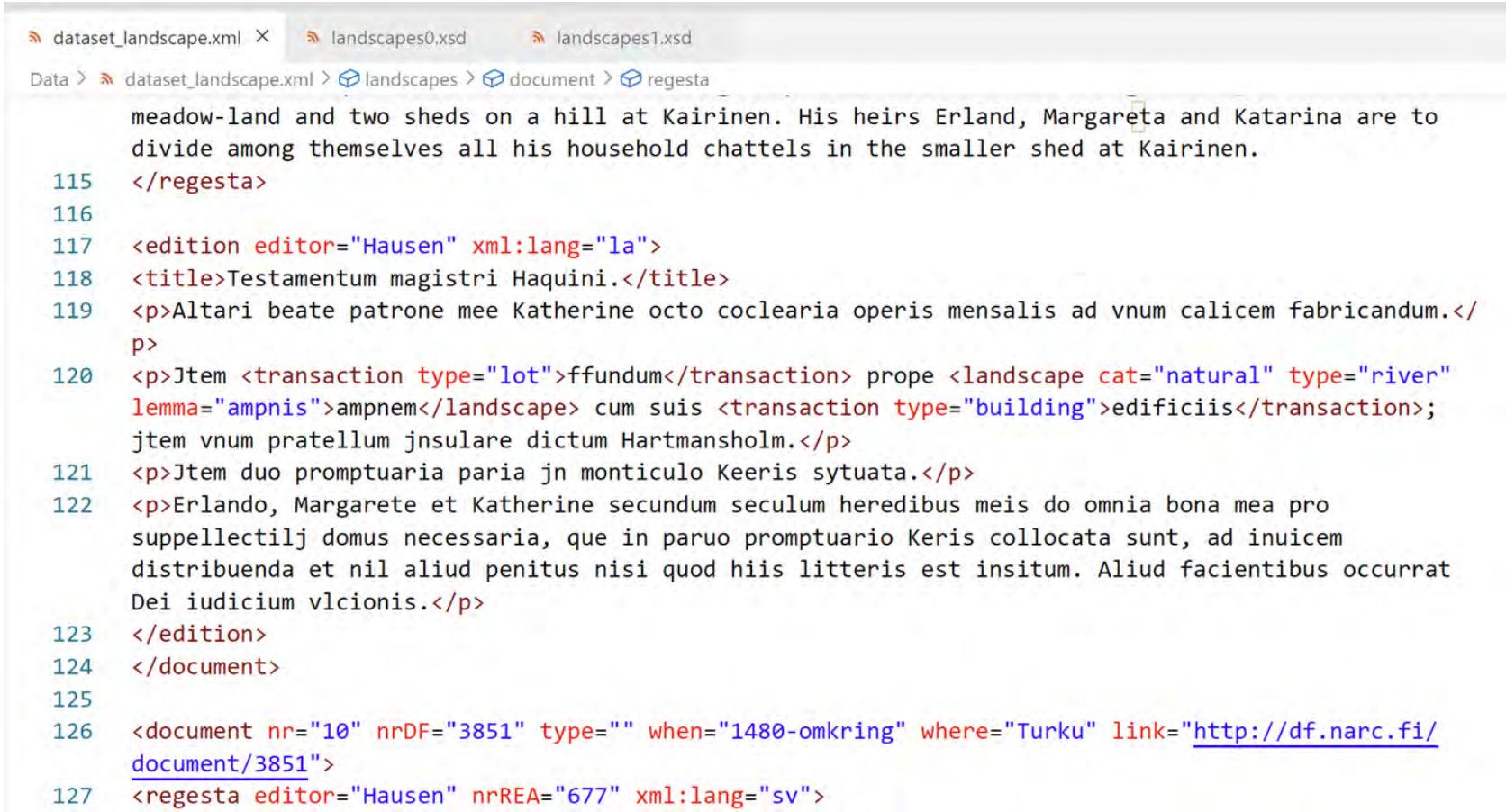


Kuvalähde: <https://aamuset.fi/artikkeli/5725211>

Steps of the Medieval Turku landmarks corpus



<https://github.com/HannaKoo/landscapes>



The screenshot shows an XML editor interface with the following details:

- File tabs: dataset_landscape.xml (active), landscapes0.xsd, landscapes1.xsd.
- Breadcrumb navigation: Data > dataset_landscape.xml > landscapes > document > regesta.
- Content pane:

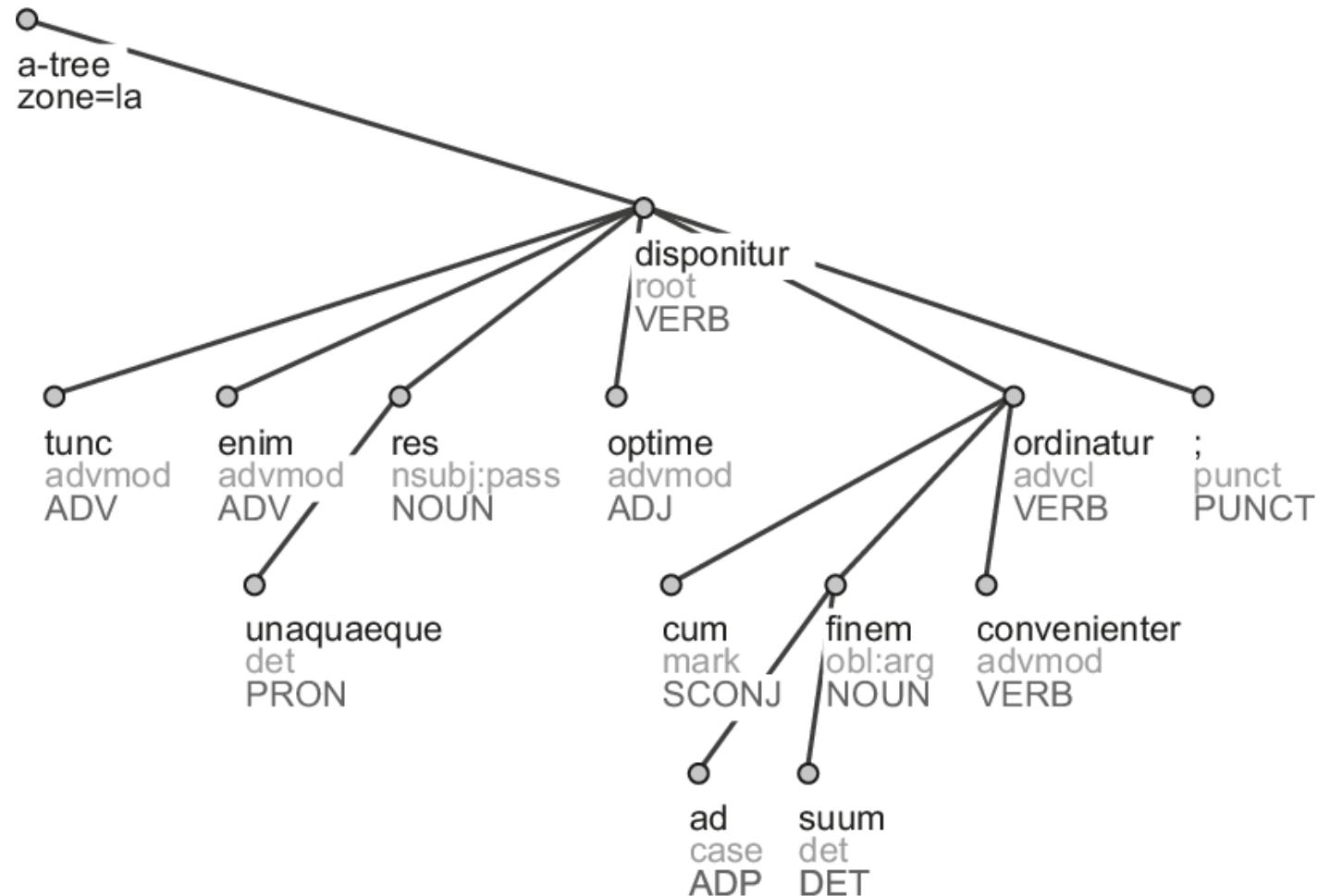
```
meadow-land and two sheds on a hill at Kairinen. His heirs Erland, Margareta and Katarina are to
divide among themselves all his household chattels in the smaller shed at Kairinen.

115 </regesta>
116
117 <edition editor="Hausen" xml:lang="la">
118 <title>Testamentum magistri Haquini.</title>
119 <p>Altari beate patronae mee Katherine octo coclearia operis mensalis ad vnum calicem fabricandum.</
p>
120 <p>Jtem <transaction type="lot">ffundum</transaction> prope <landscape cat="natural" type="river"
lemma="ampnis">ampnem</landscape> cum suis <transaction type="building">edificiis</transaction>;
jtem vnum pratellum jnsulare dictum Hartmansholm.</p>
121 <p>Jtem duo promptuaria paria jn monticulo Keeris sytuata.</p>
122 <p>Erlando, Margarete et Katherine secundum seculum heredibus meis do omnia bona mea pro
suppellectilj domus necessaria, que in paruo promptuario Keris collocata sunt, ad inuicem
distribuenda et nil aliud penitus nisi quod hiis litteris est insitum. Aliud facientibus occurrat
Dei iudicium vlcionis.</p>
123 </edition>
124 </document>
125
126 <document nr="10" nrDF="3851" type="" when="1480-omkring" where="Turku" link="http://df.narc.fi/
document/3851">
127 <regesta editor="Hausen" nrREA="677" xml:lang="sv">
```

Section 3

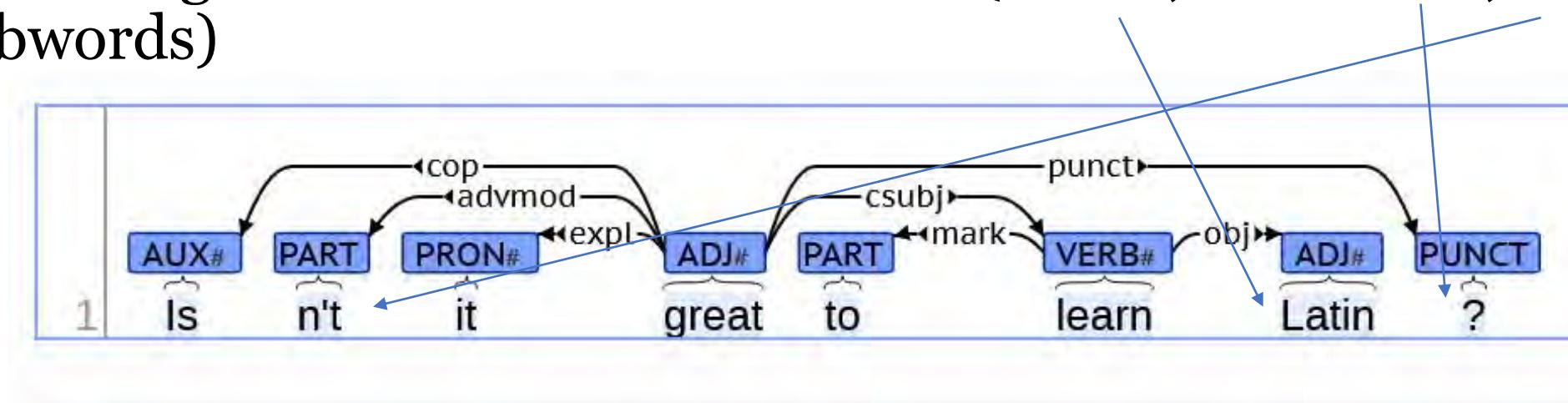
The automatic morpho-syntactic parsing

The automatic morpho-syntactic parsing



Automatic formo-syntactic parsing is:

- “to analyze the input sentence in terms of grammatical constituents, identifying the parts of speech and their syntactic relations”
- producing tokens: basic blocks of text (words, characters, or subwords)



http://epsilon-it.utu.fi/parser_demo/

Warse.org & nlp.stanford.edu

Finnish: Online parser demo

Online Parser Demo

[Turku NLP Group]

This is a demo of the Finnish dependency parsing pipeline. Pipeline includes text segmentation, morphological analysis and dependency parsing.

This demo is meant for testing only, if you need to parse a lot of text, please download the parser and model and run it locally. If you have any questions or problems contact TurkuNLP Group (contact information [here](#)).

Easiest way to run the parser locally on Linux/Mac/Windows is Docker image. See instructions [here](#).

Finnish-neural

Parser: <http://turkunlp.github.io/Turku-neural-parser-pipeline/>

Treebank documentation: <http://universaldependencies.github.io/docs/>

Enter text

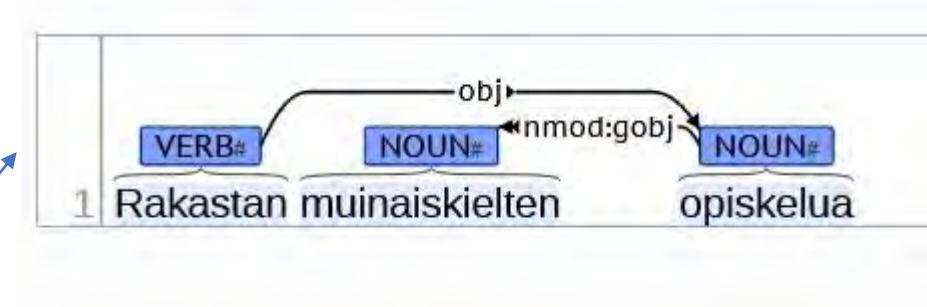
Rakastan muinaiskielten opiskelua

Parse!

Alternatively, you can upload a text file to be parsed. The results are given in a [CoNLL-U](#) formatted file.

Browse... No file selected.

Parse!



More information on the parser's project page:
<https://turkunlp.github.io/Turku-neural-parser-pipeline/>

Why use parsers?

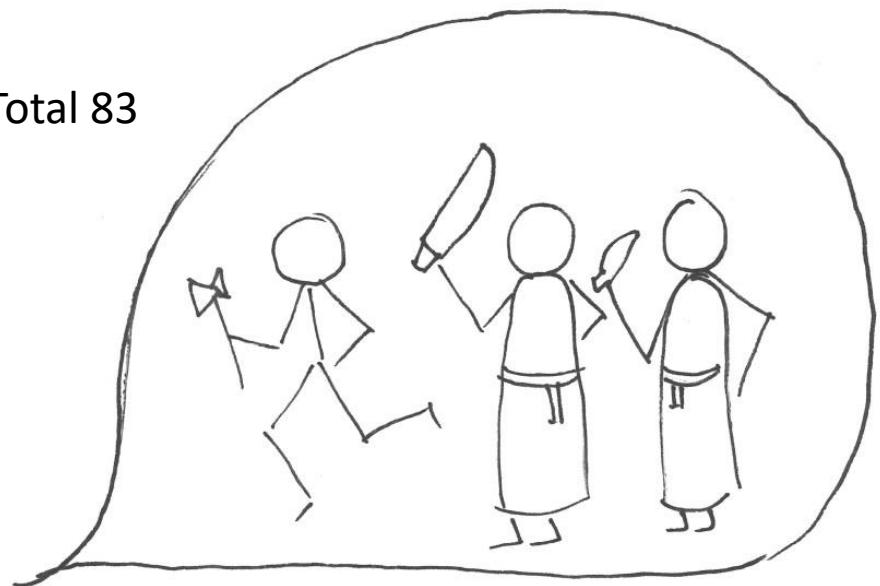


percusserit	11	119	se deffendendo aliter fugere sive evadere non valens percusserit
percusserunt	2	232	vim vi repellendo ipsum laicum adeo percusserunt
percussiendo	1	260	arrepto percussiendo in terram prostravit
percussio	1	239	quod levis fuerit percussio
percussione	23	300a	non ex percussione exponentis mors subsecuta fuerit
percussionem	1	239	prout ante percusionem predictam
percussioni	2	318	Petrus volens percussioni et periculo huiusmodi obviare
percussionibus	15	341	dictus diaconus de dictis percussionibus plene convaluerit
percussionis	1	239	ad septenam diem a die percussionis
percussis	1	239	Exinde in spatulis percussis prout ... incolumis
percussisset	2	264	quadam lancea seu cuspide percussisset
percussit	68	264	in capite ipsius percussit unico ictu
percuso	2	340	ipseque exponentis cum percuso concordaverit
percussor	1	19	Quod cum percussor suus ... viderent
percussorem	1	342	tam graviter vulneratum laicum percussorem
percussum	2	303	exponentis dictum Laurentium sic percussum dimisit
percussus	12	340	Cum autem dictus presbyter percussus
percutere	9	360	oratorem suo gladio in capite percutere voluit
percuteret	4	160	ne ipse familiaris dictum laicum percuteret
percuteretur	1	444	quodam malleo ferreo graviter percuteretur
percuti	1	300a	percuti seu forsitan interfici illum
percutiebant	1	406	pungebant et percutiebant
percutiendo	1	427	ipsum leniter percutiendo existimans
percutiens	1	360	oratorem dorsotenus percutiens
percuttere	2	145	cum deteriori parte securis percuttere voluit
percutteret	2	234	ut eum percutteret
repercussit	1	195	vim vi repellendo cum eadem securi dictum laicum taliter repercussit
repercussiens	1	292	cum cribro repercussiens et ulterius manus eius violentas dicti ... iniecit

VERBS

percusserit	11	119	se defendendo aliter fugere sive evadere non valens percusserit
percusserunt	2	232	vim vi repellendo ipsum laicum adeo percusserunt
percussiendo	1	260	arrepto percussiendo in terram prostravit
percussio	1	239	quod levis fuerit percussio
percussione	23	300a	non ex percusione exponentis mors subsecuta fuerit
percussionem	1	239	prout ante percusionem predictam
percussioni	2	318	Petrus volens percussioni et periculo huiusmodi obviare
percussionibus	15	341	dictus diaconus de dictis percussionibus plene convaluerit
percussionis	1	239	ad septenam diem a die percussionis
percussis	1	239	Exinde in spatulis percussis prout ... incolumis
percussisset	2	264	quadam lancea seu cuspidie percussisset
percussit	68	264	in capite ipsius percussit unico ictu

percutere (*percutiō*, *percussī*, *percussum*) 3 - Total 83



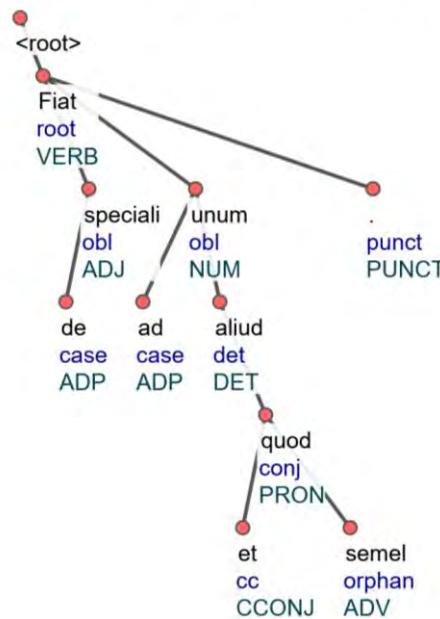
Why we need lemmas:

Conjugation of <i>dō</i> (first conjugation, irregular short <i>a</i> in most forms except <i>dās</i> and <i>dā</i>) less ▲						
indicative		singular			plural	
		first	second	third	first	second
active	present	<i>dō</i>	<i>dās</i>	<i>dat</i>	<i>damus</i>	<i>datis</i>
	imperfect	<i>dabam</i>	<i>dabās</i>	<i>dabat</i>	<i>dabāmus</i>	<i>dabātis</i>
	future	<i>dabō</i>	<i>dabis</i>	<i>dabit</i>	<i>dabimus</i>	<i>dabitis</i>
	perfect	<i>dedi</i>	<i>dedistī</i>	<i>dedit</i>	<i>dedimus</i>	<i>dedistis</i>
	pluperfect	<i>dederam</i>	<i>dederās</i>	<i>dederat</i>	<i>dederāmus</i>	<i>dederātis</i>
	future perfect	<i>dederō</i>	<i>dederis</i>	<i>dederit</i>	<i>dederimus</i>	<i>dederint</i>
passive	present	<i>dor</i>	<i>daris,</i> <i>dare</i>	<i>datur</i>	<i>damur</i>	<i>damini</i>
	imperfect	<i>dabar</i>	<i>dabāris,</i> <i>dabāre</i>	<i>dabātur</i>	<i>dabāmur</i>	<i>dabāmini</i>
	future	<i>dabor</i>	<i>daberis,</i> <i>dabere</i>	<i>dabitur</i>	<i>dabimur</i>	<i>dabimini</i>
	perfect			<i>datus</i> + present active indicative of <i>sum</i>		
	pluperfect			<i>datus</i> + imperfect active indicative of <i>sum</i>		
	future perfect			<i>datus</i> + future active indicative of <i>sum</i>		
subjunctive		singular			plural	
		first	second	third	first	second
active	present	<i>dēs,</i> <i>dūm</i>	<i>dūs,</i> <i>dūs</i>	<i>det,</i> <i>dūt</i>	<i>dēmus</i>	<i>dētis</i>
	imperfect	<i>darem</i>	<i>dārēs</i>	<i>daret</i>	<i>dārēmus</i>	<i>dārētis</i>
	perfect	<i>dederim</i>	<i>dērēs</i>	<i>dederit</i>	<i>dērēmus</i>	<i>dērētis</i>
	pluperfect	<i>dedissem</i>	<i>dēsēs</i>	<i>dedisset</i>	<i>dēsēmus</i>	<i>dēsētis</i>
passive	present	<i>der</i>	<i>dēris,</i> <i>dēre</i>	<i>dētūr</i>	<i>dēmūr</i>	<i>dēmī</i>
	imperfect	<i>darer</i>	<i>dārēs,</i> <i>dārēre</i>	<i>dārētūr</i>	<i>dārēmūr</i>	<i>dārēmī</i>
	perfect			<i>datus</i> + present active subjunctive of <i>sum</i>		
	pluperfect			<i>datus</i> + imperfect active subjunctive of <i>sum</i>		
imperative		singular			plural	
		first	second	third	first	second
active	present	—	<i>dā</i>	—	—	<i>date</i>
	future	—	<i>datō</i>	<i>datō</i>	—	<i>datōte</i>
passive	present	—	<i>dare</i>	—	—	<i>dāmī</i>
	future	—	<i>dator</i>	<i>dator</i>	—	<i>dantor</i>
non-finite forms		active			passive	
		present	perfect	future	present	perfect
infinitives		<i>dāre</i>	<i>dēisse</i>	<i>dātūrum ēsse</i>	<i>dāi</i>	<i>dātūrum ēsse</i>
	participles	<i>dāns</i>	—	<i>dātūrus</i>	—	<i>dātūs</i>
verbal nouns		gerund			supine	
		genitive	dative	accusative	ablative	accusative
		<i>dandī</i>	<i>dandō</i>	<i>dandum</i>	<i>dandō</i>	<i>datum</i>
						<i>datū</i>

Universal Dependencies

- framework for consistent annotation of grammar (parts of speech, morphological features, and syntactic dependencies) across different human languages
- open community effort:
 - 300 contributors
 - producing nearly 200 treebanks
 - over 100 languages

Fiat de speciali ad unum aliud et quod semel .



UD Treebanks for Latin

The screenshot shows a web browser window comparing five UD treebanks for Latin: UD_Latin-ITTB, UD_Latin-LLCT, UD_Latin-PROIEL, UD_Latin-Perseus, and UD_Latin-UDante. The page is titled 'treebanks/la-comparison.html' and includes a green banner stating 'This page pertains to UD version 2.' Below the banner, each treebank has a section for 'Tokenization and Word Segmentation' containing bullet points about its corpus statistics and tokenization characteristics.

UD_Latin-ITTB	UD_Latin-LLCT	UD_Latin-PROIEL	UD_Latin-Perseus	UD_Latin-UDante
Tokenization and Word Segmentation	Tokenization and Word Segmentation	Tokenization and Word Segmentation	Tokenization and Word Segmentation	Tokenization and Word Segmentation
<ul style="list-style-type: none">This corpus contains 26977 sentences and 450515 tokens.This corpus contains 62358 tokens (14%) that are not followed by a space.This corpus does not contain words with spaces.	<ul style="list-style-type: none">This corpus contains 9023 sentences, 242410 tokens and 242411 syntactic words.This corpus contains 30826 tokens (13%) that are not followed by a space.This corpus does not contain words with spaces.	<ul style="list-style-type: none">This corpus contains 18411 sentences and 200163 tokens.All tokens in this corpus are followed by a space.This corpus does not contain words with spaces.	<ul style="list-style-type: none">This corpus contains 2273 sentences and 29138 tokens.This corpus contains 4381 tokens (15%) that are not followed by a space.This corpus does not contain words with spaces.	<ul style="list-style-type: none">This corpus contains 1721 sentences, 55287 tokens and 55524 syntactic words.This corpus contains 8537 tokens (15%) that are not followed by a space.This corpus does not contain words with spaces.

Simple example:

Puella canes amat.
The girl loves dogs.

The screenshot shows the UD Pipe service page on the LINDAT/CLARIN website. The URL in the browser is https://lindat.mff.cuni.cz/services/udpipe/. The page has a dark blue header with the LINDAT logo, search, catalogue, education, projects, tools, services, and about links. It also features Dariah-EU and CLARIN logos. The main content area has a light gray background. At the top, it says "UDPipe". Below that are three buttons: "About", "Run" (which is selected), and "REST API Documentation". A text block explains UD Pipe is a trainable pipeline for tokenization, tagging, lemmatization and dependency parsing of CoNLL-U files. It mentions it's language-agnostic and can be trained given annotated data in CoNLL-U format. Trained models are provided for nearly all UD treebanks. UD Pipe is available as a binary for Linux/Windows/OS X, as a library for C++, Python, Perl, Java, C#, and as a web service. A third-party R CRAN package also exists. Another text block states UD Pipe is free software distributed under the Mozilla Public License 2.0 and the linguistic models are free for non-commercial use and distributed under the CC BY-NC-SA license, although for some models the original data used to create the model may impose additional licensing conditions. UD Pipe is versioned using Semantic Versioning. Copyright 2017 by Institute of Formal and Applied Linguistics, Faculty of Mathematics and Physics, Charles University, Czech Republic. A note says the service is freely available for testing and respects the CC BY-NC-SA licence of the models – explicit written permission of the authors is required for any commercial exploitation of the system. If you use the service, you agree that data obtained by us during such use can be used for further improvements of the systems at UFAL. All comments and reactions are welcome. Below this, there are options to select a model (UD 2.10, UD 2.6, EvaLatin20) and an input file (latin-proiel-ud-2.10-220711). Under actions, checkboxes are checked for "Tag and Lemmatize" and "Parse". A "Advanced Options" button is at the bottom.

CoNLL-U output

<https://lindat.mff.cuni.cz/services/udpipe/>

```
# generator = UDPipe 2, https://lindat.mff.cuni.cz/services/udpipe
# udpipe_model = latin-perseus-ud-2.10-220711
```

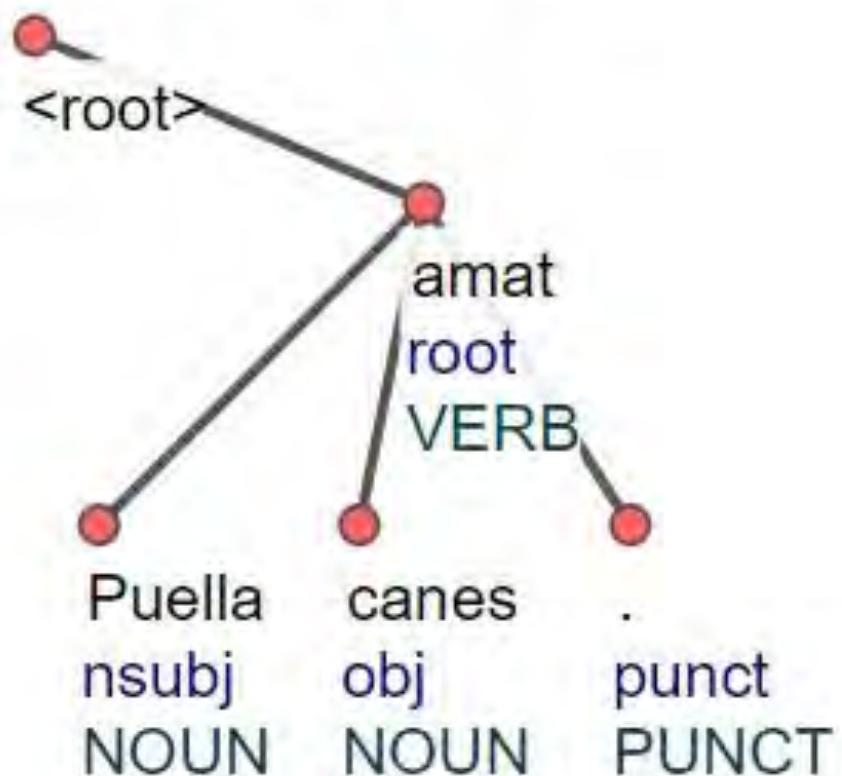
```
# sent_id = 1
# text = Puella canes amat.
1 Puella puella NOUN
2 canes canis NOUN
2 amat amo VERB
3 4 . . PUNCT
```

The screenshot shows the UDPIPE service interface at <https://lindat.mff.cuni.cz/services/udpipe/>. The top navigation bar includes links for UD 2.10, UD 2.6, and EvaLatin20, along with a search bar and a magnifying glass icon. The main area is titled 'Service' and contains a message about respecting the CC BY-NC-SA license and giving explicit written permission for commercial exploitation. Below this, the 'Model' dropdown is set to 'UD 2.10 (description)', and the 'Actions' checkboxes for 'Tag and Lemmatize' and 'Parse' are checked. The 'Input Text' field contains the Latin sentence 'Puella canes amat.'. A green button labeled 'Process Input' is visible above the output area.



Tree

Puella canes amat .



Lemmatization and morphological analysis for the Latin Dependency Treebank (Celano)

- lemmatization can be defined as consisting in the assignment of an ‘ID word form’ – the dictionary head word
- **dō** (*present infinitive dare, perfect active dedi, supine datum*); *1st conjugation* – you can **not** truncate your search at d*
- Dictionary head word **amo**

Ōdī et amō. Quārē id faciam fortasse requīris.
Nesciō, sed fierī sentiō et excrucior. (Catullus)

Participles:

- The annotation of participles is demanding: *particeps* 'to take part in the nature of both verb and noun'
- participle **amāns** -> **amo** VERB 'loving'
- **amāns** m/f III -> NOUN 'lover'
- ...sed nihil difficile **amanti** puto 'I believe nothing to be difficult for the dedicated' Cicero (*Orator ad M. Brutum*)
- And when to consider it an adjective? - > ADJ
- Maria, femina **laundanda**, advenit. 'Maria, a praiseworthy woman, has arrived.'

Discussion

- Lemmatization and morphological analysis depends on digitized dictionaries (Lewis and Short, 1879)
- Printed dictionaries rely on providing definitions for humans (Preferably with holistic understanding)
- Treebanks need to be consistent
- Latin dependency treebank relies on rule based tokenization and sentence split algorithms, whose output feeds the COMBO lemmatizer PoS tagger and parser

Ensemble lemmatization with the Classical Language Toolkit (Burns)

- A combination of data-driven and rule based approaches
 - Ensemble lemmatizer
 - How to handle unseen vocabulary?
 - Especially for Latin: resources are extremely limited for dialect, period and genre
 - There is no effective way to combine tools
 - Backoff lemmatization reflects the methods of the philologically educated reader of historical languages

Creating a Gold Standard

- Circa 1 000 words
- 500 chosen on chance
- 500 on deliberately picked
- Lemmas – easy
- POS – easy
- Morphology – a bit of thinking
- Syntax – guaranteed challenge



UDante and UD pipe

- 134/1 251 lemma
-> 10 % needed to be corrected
 - 63/1 251 POS -> 5 %

Key take home messages:

- We want to ensure as much consistency as possible *id est*:
 - human use allows for a number of irregularities that impinge on computational resources derived from them
 - There are no guidelines for any of the UD treebanks in Latin – a new global effort – UD4HL

Future outlooks

- MDA Multi-Dimensional Analysis -> registers
- Social history and everyday life in medieval Turku
- Linguistics landscaped data on Kielipankki and opportunities for new research topics

Corpus of landscapes in medieval documents from Turku, source

View resource name in all available languages

medievalturku

Persistent Identifier of this resource: <http://urn.fi/urn:nbn:fi:lb-2023032021>

This resource will be made available in Kielipankki – the Language Bank of Finland.

This is a text corpus for the study of the linguistic landscapes of Turku. The material is mostly based on the so-called Black Book of the Turku Cathedral, but also on other contemporary documents concerning the... [Read More](#)

View resource description in all available languages

[« Back](#) [Edit Resource](#)

text	
Distribution	Multilingual text corpus
Availability	Languages
Under Negotiation	Latin
License	Swedish
CLARIN PUB	Variety: medieval written Swedish (Type: Other) (20,000 Words)
Licensor	Linguality
University of Turku	Linguality type: Multilingual Multi-linguality type: Other
IPR Holder	Size
University of Turku	86 Texts
Contact Person	
Hanna-Mari Kupari	

Resource Creation

Resource Creator

Hanna-Mari Kupari
Marko Lamberg

Metadata

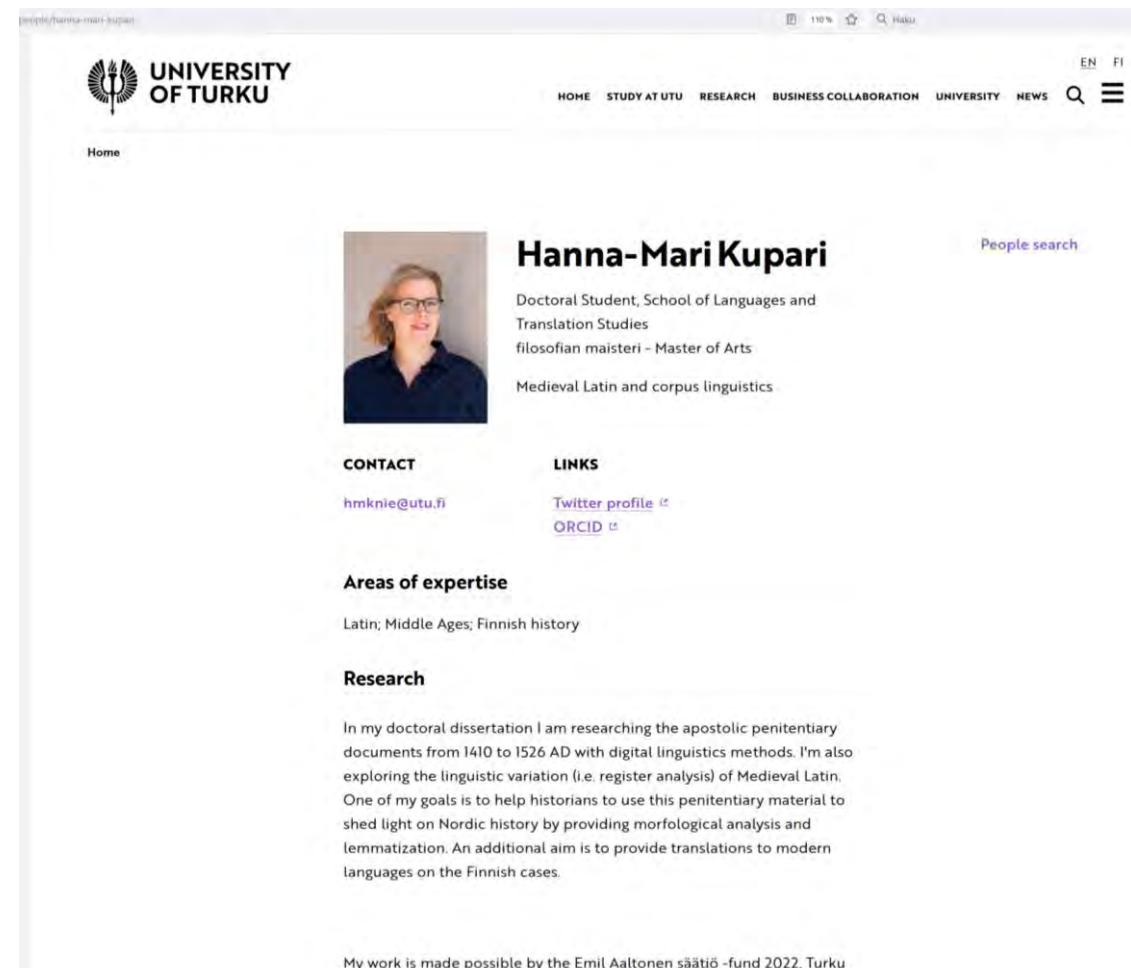
Created: 03/20/2023
Last Updated: 03/22/2023

Metadata Creator

Mietta Lennes

Thanks

- hmknie@utu.fi
- <https://github.com/HannaKoo>
- Twitter: @kuparimari
- Instagram:
[hannasuperturkulainen](https://www.instagram.com/hannasuperturkulainen)
- Kalmistopiiri



The image shows a screenshot of the University of Turku website. The header features the university's logo and navigation links for Home, Study at UTU, Research, Business Collaboration, University, News, and a search bar. The main content area displays a profile for Hanna-Mari Kupari, showing her photo, name, title (Doctoral Student, School of Languages and Translation Studies), and research interests (Medieval Latin and corpus linguistics). It also includes sections for Contact (email: hmknie@utu.fi) and Links (Twitter profile and ORCID). Below this, there are sections for Areas of expertise (Latin; Middle Ages; Finnish history) and Research, which describes her work on apostolic penitentiary documents and Medieval Latin.

people/hanna-mari-kupari

UNIVERSITY OF TURKU

Home

Hanna-Mari Kupari

Doctoral Student, School of Languages and Translation Studies
filosofian maisteri - Master of Arts
Medieval Latin and corpus linguistics

CONTACT

hmknie@utu.fi

LINKS

[Twitter profile](#) (2)
[ORCID](#) (2)

Areas of expertise

Latin; Middle Ages; Finnish history

Research

In my doctoral dissertation I am researching the apostolic penitentiary documents from 1410 to 1526 AD with digital linguistics methods. I'm also exploring the linguistic variation (i.e. register analysis) of Medieval Latin. One of my goals is to help historians to use this penitentiary material to shed light on Nordic history by providing morphological analysis and lemmatization. An additional aim is to provide translations to modern languages on the Finnish cases.

My work is made possible by the Emil Aaltonen säätiö -fund 2022, Turku