



UNIVERSITÀ  
CATTOLICA  
del Sacro Cuore



# Latin Lemmatization & POS Tagging

Issues, Resources, Tools

Francesco Mambrini, Marco Passarotti

*IANLS Vacation School: Digital Humanities and Neo-Latin Studies*  
University of Bonn



This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme - Grant Agreement No. 769994.

## Lemmatization & Part-of-Speech Tagging

### Textual Resources

Lemmatized Corpora for Latin

### Tools and Hands-on

Tools for lemmatization & POS Tagging

Try it yourself on the provided text(s)!

### Latin in the Semantic Web

The LiLa Knowledge Base

The TextLinker

### Querying Lemmatized Resources

Latin Treebanks in Universal Dependencies

Interlinked Lexical and Textual Resources in LiLa

## Goals

**Lemmatization** and **part-of-speech tagging** (POS-tagging) aim to **abstract** some linguistic properties to allow **form-invariant** reference to types/tokens.

- ?! How can I retrieve all the occurrences of a word in a text?
- ?! How can I know which (morphosyntactic) function(s) a word plays in a text?

Different word forms in different contexts...

- ▶ ... *his rebus cognitis Caesar Gallorum animos **verbis** confirmavit...*  
→ ablative plural (token); dative & ablative plural (type)
- ▶ ... *quod ego si **verbo** adsequi possem...*  
→ ablative singular (token); dative & ablative singular (type)
- ▶ ... *ne more iuvencae mugiat, et timide **verba** intermissa retemptat...*  
→ accusative plural (token); nominative, accusative & vocative plural (type)

...but all can be referred to a canonical/standardized citation form  
(Lemma):

⇒ **uerbum**

→ nominative singular of neuter II. declension noun

What about *cognitis*, *intermissa* and *timide*?

## Lemmatization

**Type-based** the process of assigning each type (in a text) to one, or more lemma(s)

**Token-based** the process of assigning each token in a text to a lemma

Different lexicographic criteria:

- ▶ inflectional morphology: same paradigm, same lemma? what about participles?
- ▶ graphical representation: *voluptas* vs. *uoluptas*
- ▶ spelling: *sulphur* vs. *sulfur*
- ▶ ending and inflectional type: *diameter* vs. *diametros* vs. *diametrus*
- ▶ paradigmatic slot for the lemma: *sequor* vs. *sequo* (see Du Cange: infinitives used)
- ▶ homographs: *occido*/[*caedo*|*cado*] vs. *occido*[1|2]

Words can play different (morphosyntactic) functions in sentences:

★ **supra**

- ▶ ...*ager trecentis aut etiam **supra** nummorum milibus emptus...*  
→ adverb (ADV)
- ▶ ...*ille qui **supra** nos habitat...*  
→ preposition (ADP)

★ **scribo**

- ▶ ...*atque in Thesauro **scripsit** causam dicere prius unde petitur...*  
→ verb (VERB)

★ **elephantus**

- ▶ ...***elephanto** beluarum nulla prudentior...*  
→ noun (NOUN)

These functions are predictable and come from a rather small set of alternatives.

## Part-of-speech tagging

**Type-based:** the process of assigning each type one, or more morphosyntactic **function(s)**, i. e. parts of speech, from a given set

**Token-based** the process of assigning each token in a text one morphosyntactic **function**, i. e. part of speech, from a given set

Current standard de facto tagset: **Universal Dependencies**



16+1 classes: ADJ (*adjectives*), ADP (*pre- & postpositions*), ADV (*adverbs*), AUX (*auxiliaries*), CCONJ & SCONJ (*co-ordinating & subordinating conjunctions*), DET (*determiners*), INTJ (*interjections*), NOUN & PROPN (*common & proper nouns*), NUM (*numerals*), PART (*particles*), PRON (*pronouns*), VERB (*verbs*), SYM (*symbols*), X (*other*) + PUNCT (*punctuation*)

<https://universaldependencies.org>

Type-based vs. token-based POS tagging:

- ▶ Every ADJ can be NOUN
- ▶ Every ADP, CCONJ, SCONJ etc. can be NOUN (like in metalinguistic discourse)
- ▶ Every VERB can be NOUN

One or more part-of-speech? Which part-of-speech?

- ▶ *italicus*: ADJ? NOUN? PROPEN?
- ▶ *ubi*: ADV? SCONJ?
- ▶ *non*: ADV? PART?
- ▶ *aliqui*: PRON? DET? ADJ?

Lemmatization & Part-of-Speech Tagging

## Textual Resources

Lemmatized Corpora for Latin

## Tools and Hands-on

Tools for lemmatization & POS Tagging

Try it yourself on the provided text(s)!

## Latin in the Semantic Web

The LiLa Knowledge Base

The TextLinker

## Querying Lemmatized Resources

Latin Treebanks in Universal Dependencies

Interlinked Lexical and Textual Resources in LiLa

- ▶ LASLA Corpus
- ▶ Index Thomisticus
- ▶ Computational Historical Semantics
- ▶ 5 Latin Treebanks in UD
- ▶ CLaSSES
- ▶ ...and others

Lemmatization & Part-of-Speech Tagging

Textual Resources

Lemmatized Corpora for Latin

**Tools and Hands-on**

Tools for lemmatization & POS Tagging

Try it yourself on the provided text(s)!

Latin in the Semantic Web

The LiLa Knowledge Base

The TextLinker

Querying Lemmatized Resources

Latin Treebanks in Universal Dependencies

Interlinked Lexical and Textual Resources in LiLa

Analysed wordform : sulphur

```
=====ANALYSIS =====
SEGMENTATION:  sulphur

-----morphological feats 1 -----
--nns--

Case:  Nominative
Gender: Neuter
Number: Singular
-----morphological feats 2 -----
--ans--

Case:  Accusative
Gender: Neuter
Number: Singular
-----morphological feats 3 -----
--vns--

Case:  Vocative
Gender: Neuter
Number: Singular
=====LEMMMA =====
sulpur                N3B  s3429 n
-----morphological feats-----
NcC

PoS:  Noun
Type:  Common
Inflexional Category:  III decl
```

```
# generator = UDPipe 2, https://lindat.mff.cuni.cz/services/udpipe
# udpipe_model = latin-proiel-ud-2.6-200830
# udpipe_model_licence = CC BY-NC-SA
# newdoc
# newpar
# sent_id = 1
# text = Cui dono lepidum novum libellum arida modo pumice expol
1  Cui qui PRON Pr Case=Dat|Gender=Masc|Number=Sing
2  dono donum NOUN Nb Case=Abl|Gender=Neut|Number=Neut
3  lepidum lepidus ADJ A- Case=Acc|Degree=Pos|Gender=Masc|Number=Neut
4  novum novus ADJ A- Case=Acc|Degree=Pos|Gender=Masc|Number=Neut
5  libellum libellus NOUN Nb Case=Acc|Gender=Masc|Number=Neut
6  arida aridus ADJ A- Case=Acc|Degree=Pos|Gender=Masc|Number=Neut
7  modo modo ADV Df _ 8 advmod _ TokenFrequency=1
8  pumice pumic NOUN Nb Case=Abl|Gender=Masc|Number=Neut
9  expolitur? expolio VERB V- Case=Nom|Gender=Masc|Number=Neut
SpaceAfter=No|TokenRange=51:61
```

```
# generator = UDPipe 2, https://lindat.mff.cuni.cz/services/udpipe
# udpipes_model = latin-evalatin20-200830
# udpipes_model_licence = CC BY-NC-SA
# newdoc
# newpar
# sent_id = 1
# text = Cui dono lepidum novum libellum arida modo pumice expolitum?
1  Cui qui PRON _ _ _ _ _ TokenRange=0:3
2  dono donum NOUN _ _ _ _ _ TokenRange=4:8
3  lepidum lepidus ADJ _ _ _ _ _ TokenRange=9:16
4  novum novus ADJ _ _ _ _ _ TokenRange=17:22
5  libellum libellus NOUN _ _ _ _ _ SpacesAfter=\r\n|TokenRange=23:31
6  arida aridus ADJ _ _ _ _ _ TokenRange=33:38
7  modo modo ADV _ _ _ _ _ TokenRange=39:43
8  pumice pumicus NOUN _ _ _ _ _ TokenRange=44:50
9  expolitum? expolito VERB _ _ _ _ _ SpaceAfter=No|TokenRange=51:61
```

- ▶ Download the tool from <https://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>
- ▶ Prepare a txt file with a Latin text
- ▶ Tokenize the file and prepare the input (one-word-per-line):  

```
cd treetagger/cmd  
perl utf8-tokenize.perl INPUT-FILE.txt >  
OUTPUT-FILE.txt
```

▶ `cd ../bin`

▶ Linux/Mac:

```
./tree-tagger <parameter-file> <input-file>  
<output-file> -token -lemma
```

Example (download a parameter file for Latin and put it into the 'bin' folder): `./tree-tagger latin.par input.txt output.txt -token -lemma`

▶ Windows:

```
tag-LANGUAGE.bat <input-file> <output-file>  
Example: tag-latin.bat input.txt output.txt
```

- ▶ Collatinus Web:  
<https://outils.biblissima.fr/en/collatinus-web/>
- ▶ Deucalion: <https://dh.chartes.psl.eu/deucalion/latin>
- ▶ Stanza: three models for Latin. [https://stanfordnlp.github.io/stanza/available\\_models.html](https://stanfordnlp.github.io/stanza/available_models.html)
- ▶ Morpheus: <https://github.com/PerseusDL/morpheus>
- ▶ Whitaker's Words: <https://latin-words.com>

Lemmatization & Part-of-Speech Tagging

Textual Resources

Lemmatized Corpora for Latin

**Tools and Hands-on**

Tools for lemmatization & POS Tagging

Try it yourself on the provided text(s)!

Latin in the Semantic Web

The LiLa Knowledge Base

The TextLinker

Querying Lemmatized Resources

Latin Treebanks in Universal Dependencies

Interlinked Lexical and Textual Resources in LiLa

## Lemmatization & Part-of-Speech Tagging

### Textual Resources

Lemmatized Corpora for Latin

### Tools and Hands-on

Tools for lemmatization & POS Tagging

Try it yourself on the provided text(s)!

## Latin in the Semantic Web

The LiLa Knowledge Base

The TextLinker

### Querying Lemmatized Resources

Latin Treebanks in Universal Dependencies

Interlinked Lexical and Textual Resources in LiLa

# The Linked Data Principles

...just to be FAIR



- ▶ Use URIs for things (e.g. an entry in a lexicon, a token in a corpus)

# The Linked Data Principles

...just to be FAIR



- ▶ Use URIs for things (e.g. an entry in a lexicon, a token in a corpus)
- ▶ Use HTTP URIs to allow people (and machines) to look up things

# The Linked Data Principles

...just to be FAIR



- ▶ Use URIs for things (e.g. an entry in a lexicon, a token in a corpus)
- ▶ Use HTTP URIs to allow people (and machines) to look up things
- ▶ Use web standards to represent/query (meta)data, such as RDF and SPARQL

# The Linked Data Principles

...just to be FAIR



- ▶ Use URIs for things (e.g. an entry in a lexicon, a token in a corpus)
- ▶ Use HTTP URIs to allow people (and machines) to look up things
- ▶ Use web standards to represent/query (meta)data, such as RDF and SPARQL
- ▶ Include links to other URIs

# Why To Apply LD to Linguistic Resources

J. Gracia: LLD CLARIN Café, 29/4/21



# Why To Apply LD to Linguistic Resources

J. Gracia: LLD CLARIN Café, 29/4/21



- ▶ Resources disconnected from each other (silos of LRs)

# Why To Apply LD to Linguistic Resources

J. Gracia: LLD CLARIN Café, 29/4/21



- ▶ Resources disconnected from each other (silos of LRs)
- ▶ Proprietary and heterogeneous formats

# Why To Apply LD to Linguistic Resources

J. Gracia: LLD CLARIN Café, 29/4/21



- ▶ Resources disconnected from each other (silos of LRs)
- ▶ Proprietary and heterogeneous formats
- ▶ Different representation schemes, query languages, annotation criteria and tagsets

# Benefits of Applying LD to Linguistic Resources

Chiarcos et al. (2013)



# Benefits of Applying LD to Linguistic Resources

Chiarcos et al. (2013)



- ▶ Representation and Modelling: RDF is a very versatile data model to represent stand-off annotations, dependency parses etc.

# Benefits of Applying LD to Linguistic Resources

Chiarcos et al. (2013)



- ▶ Representation and Modelling: RDF is a very versatile data model to represent stand-off annotations, dependency parses etc.
- ▶ Structural Interoperability: HTTP, URIs, RDF, SPARQL

# Benefits of Applying LD to Linguistic Resources

Chiarcos et al. (2013)



- ▶ Representation and Modelling: RDF is a very versatile data model to represent stand-off annotations, dependency parses etc.
- ▶ Structural Interoperability: HTTP, URIs, RDF, SPARQL
- ▶ Conceptual Interoperability: common ontologies and re-usable vocabularies to understand how to use the URIs: resources are explicitly linked

# Benefits of Applying LD to Linguistic Resources

Chiarcos et al. (2013)



- ▶ Representation and Modelling: RDF is a very versatile data model to represent stand-off annotations, dependency parses etc.
- ▶ Structural Interoperability: HTTP, URIs, RDF, SPARQL
- ▶ Conceptual Interoperability: common ontologies and re-usable vocabularies to understand how to use the URIs: resources are explicitly linked
- ▶ Federation: to combine information from physically separated repositories

# Benefits of Applying LD to Linguistic Resources

Chiarcos et al. (2013)



- ▶ Representation and Modelling: RDF is a very versatile data model to represent stand-off annotations, dependency parses etc.
- ▶ Structural Interoperability: HTTP, URIs, RDF, SPARQL
- ▶ Conceptual Interoperability: common ontologies and re-usable vocabularies to understand how to use the URIs: resources are explicitly linked
- ▶ Federation: to combine information from physically separated repositories
- ▶ Dynamicity: to provide access to the most recent version of a resource

# Benefits of Applying LD to Linguistic Resources

Chiarcos et al. (2013)



- ▶ Representation and Modelling: RDF is a very versatile data model to represent stand-off annotations, dependency parses etc.
- ▶ Structural Interoperability: HTTP, URIs, RDF, SPARQL
- ▶ Conceptual Interoperability: common ontologies and re-usable vocabularies to understand how to use the URIs: resources are explicitly linked
- ▶ Federation: to combine information from physically separated repositories
- ▶ Dynamicity: to provide access to the most recent version of a resource
- ▶ Ecosystem: a large and active community with common tools and practices. Initiatives: (1) COST Action *Nexus Linguarum* (COST Action 2019-2023): European network for Web-centred linguistic data science; (2) *Prêt-à-LLOD* (RIA 2019-2022): Ready-to-use Multilingual Linked Language Data for Knowledge Services across Sectors; (3) LD4LT (*Linked Data for Language Technology Community Group*): to create a consolidated LOD vocabulary for web (linguistic) annotation

## ERC Consolidator Grant 2018-2023

A collection of multifarious, interoperable linguistic resources described with the same vocabulary for knowledge description (by using common data categories and ontologies)

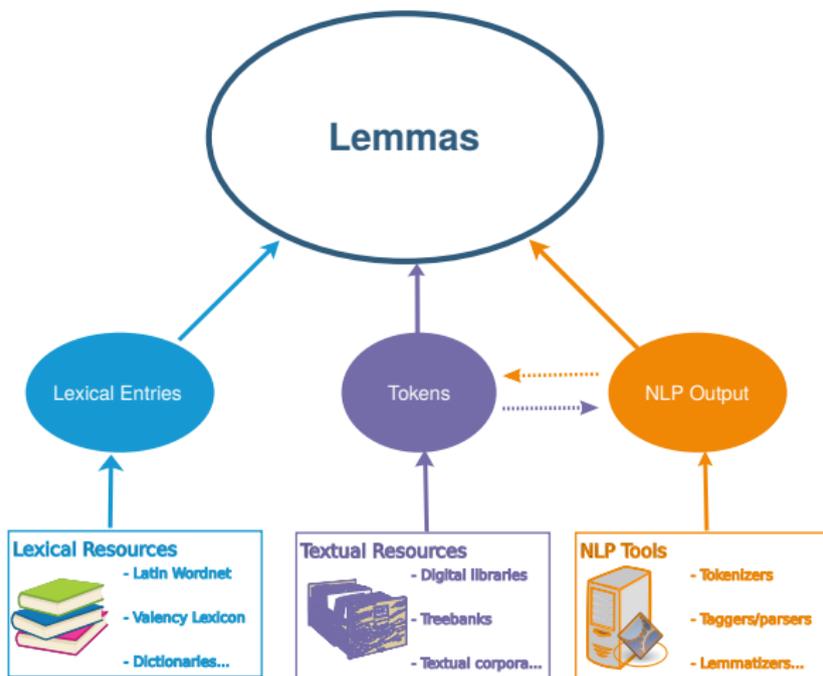
### Interlinking as a Form of Interaction

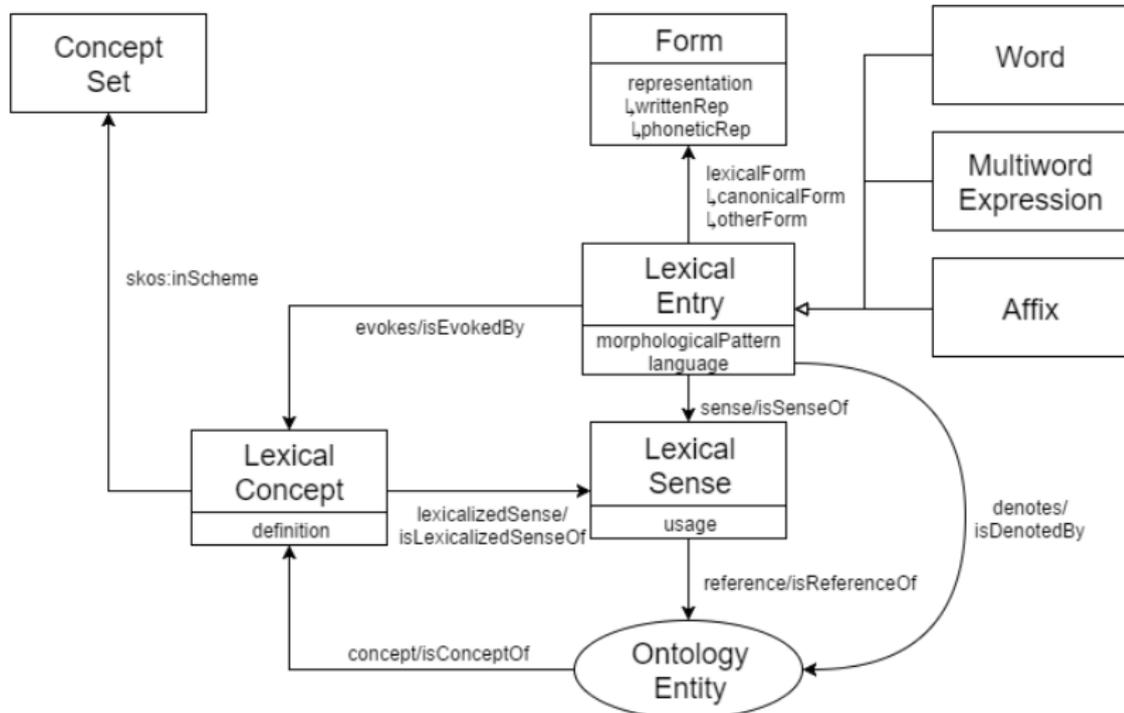


Infrastructure



Interoperability





Lemma *admiror* 'to admire, to respect'

<http://lila-erc.eu/data/id/lemma/87541>

- ▶ Lemma Bank
- ▶ A bilingual dictionary (Lewis & Short)
- ▶ A derivational lexicon (Word Formation Latin)
- ▶ A polarity lexicon (LatinAffectus)
- ▶ An etymological dictionary (De Vaan)
- ▶ A Valency Lexicon (Latin Vallex)
- ▶ A manually checked subset of the Latin WordNet

## Lemma Bank Query Interface

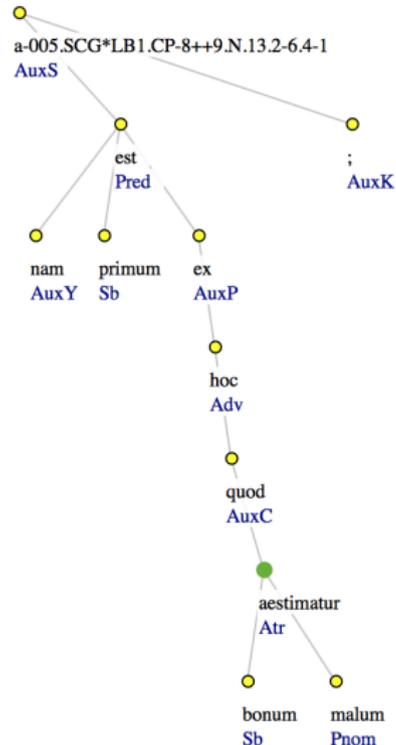
<https://lila-erc.eu/query/>

# Textual Resources

Source: the *Index Thomisticus* Treebank (original scheme)

*nam primum est ex hoc  
quod bonum aestimatur  
malum;* (IT-TB: SCG, lib. 1,  
cap. 89, n. 13)

*for the first arises because  
the good is judged to be  
evil;* (Trans. Anton C. Pegis)





## Token *aestimatur*

[http://lila-erc.eu/lodview/data/corpora/  
ITTB/id/token/005.SCG\\*LB1.CP-8++9.N.13.  
2-6.4-1W8](http://lila-erc.eu/lodview/data/corpora/ITTB/id/token/005.SCG*LB1.CP-8++9.N.13.2-6.4-1W8)

## ▶ Textual Resources

- ✓ Index Thomisticus Treebank (*Summa contra Gentiles*): ca. 400,000 nodes
- ✓ UDante Treebank: ca. 46,000 tokens
- ✓ *Querolus sive Aulularia*: ca. 17,000 tokens
- ✓ *Liber Abbaci* (ch. VIII) by Leonardo Fibonacci: ca. 30,000 tokens
- ✓ LASLA Corpus: ca. 1.7 million tokens
- PROIEL and LLCT treebanks, CompHistSem, CroALa, Musisque DeoQue

## ▶ Lexical Resources

- ✓ Lemma Bank: ca. 200,000 canonical forms
- ✓ Word Formation Latin: ca. 36,000 lemmas (Classical Latin)
- ✓ Etymological Dict. of Latin & the Other Italic Langs.: ca. 1,500 entries
- ✓ LatinAffectus: ca. 3,300 entries
- ✓ Index Graecorum Vocabulorum in L. Latinam Transl.: ca. 1,800 entries
- ✓ Latin WordNet: ca. 2,500 manually checked entries
- ✓ Latin Vallex 2.0: ca. 2,000 entries
- ✓ Lewis & Short Dictionary: ca. 50,000 entries
- Lexikon der Indogermanischen Verben (LIV). Wiktionary, BabelNet

**TOTAL: approximately 33 million triples**

## Lemmatization & Part-of-Speech Tagging

### Textual Resources

Lemmatized Corpora for Latin

### Tools and Hands-on

Tools for lemmatization & POS Tagging

Try it yourself on the provided text(s)!

## Latin in the Semantic Web

The LiLa Knowledge Base

The TextLinker

## Querying Lemmatized Resources

Latin Treebanks in Universal Dependencies

Interlinked Lexical and Textual Resources in LiLa



LILA: TEXT LINKER (β)

LiLa  
Linking Latin

PASTE YOUR TEXT BELOW

TEXT PROCESS

Vivamus mea Lesbia, atque amemus,  
rumoresque senum severiorum  
omnes unius aestimemus assis!  
soles occidere et redire possunt:  
nobis cum semel occidit brevis lux,  
nox est perpetua una dormienda.  
da mi basia mille, deinde centum,  
dein mille altera, dein secunda centum,  
deinde usque altera mille, deinde centum.  
dein, cum milia multa fecerimus,  
conturbabimus illa, ne sciamus,  
aut ne quis malus invidere possit,  
cum tantum sciat esse basiorum. |

Copyright © LiLa ERC 2020

Figure: LiLa's Text Linker



**LILA: TEXT LINKER (β)**

PASTE YOUR TEXT BELOW

TEXT PROCESS

Vivamus nea Lesbia , atque amemus , rumoresque senum severiorum omnes unius aestinemus assis !  
soles occidere et redire possunt :  
nobis cum semel occidit brevis lux , nox est perpetua una dormienda .  
da mi basia mille , deinde centum , dein mille altera , dein secunda centum , deinde usque  
altera mille , deinde centum .  
dein , cum milia multa fecerimus , conturbabimus illa , ne sciamus , aut ne quis malus  
invidere possit , cum tantum sciāt esse basiorum .

LILA KNOWLEDGE BASE LINKING

Click a token to show linked data

Form: basia

Lemma: basium - Upos: NOUN

Data from LemmaBank:

Linked to LiLa [lilaLemma:91284](#)

rdfs:type Lemma  
rdfs:label basium  
lila:hasBase Base536  
lila:hasGender neuter

Legend:  
exact match  
ambiguous match  
no match

Copyright © LiLa ERC 2020

Figure: Text processed against the LiLa Knowledge Base

Try the TextLinker yourself on the provided text(s)!

## Lemmatization & Part-of-Speech Tagging

### Textual Resources

Lemmatized Corpora for Latin

### Tools and Hands-on

Tools for lemmatization & POS Tagging

Try it yourself on the provided text(s)!

### Latin in the Semantic Web

The LiLa Knowledge Base

The TextLinker

### Querying Lemmatized Resources

Latin Treebanks in Universal Dependencies

Interlinked Lexical and Textual Resources in LiLa

Select one UD Latin treebank from the list

- ▶ `dicitur`
- ▶ `L=dico`
- ▶ `NOUN`
- ▶ `L=dico & Number=Sing`
- ▶ `L=dico &! Number=Sing`
- ▶ `L=dico | L=materia`
- ▶ `L=dico >nsubj L=commentator (ITTB)`

Select collection, language (Latin) and corpus

- ▶ Basic: worform (hominem) or, if the query is a lemma (homo), all the forms of that lemma
- ▶ Lemma: homo
- ▶ Phrase: forma de
- ▶ Word part: hom
- ▶ CQL: [upos="NOUN" & (lemma="homo" | lemma="forma")]

You can always specify the context  
(lemmas co-occurring with query results in a specified window size)

## Select language (Latin) and corpus

- ▶ `pattern {N [lemma="homo"]}`
- ▶ `pattern {N [upos="NOUN"]}`
- ▶ `pattern {V [upos=VERB];} without {V [lemma="sum"]}`
- ▶ `pattern {N1 [lemma="forma"]; N2 [lemma="materia"]; N1 < N2 }`
- ▶ `pattern {N1 [upos=VERB]; N2 [upos=NOUN]; N1 < N2 }`
- ▶ `pattern {GOV -> DEP; DEP [upos=NOUN]}`
- ▶ `pattern {N [upos = AUX]}`

Select collection (UD), language (Latin) and corpus

- ▶ a-node \$n1:= [ lemma = "homo" ]
- ▶ a-node \$n1:= [ conll/cpos = "NOUN" ]
- ▶ a-node \$n1:= [ conll/cpos = "NOUN" ]  
» for \$n1.lemma give \$1, count(), sort by \$2 desc, \$1
- ▶ a-node \$n1:= [ conll/cpos = "NOUN", child \$n2:= [ conll/depel = "amod" ] ]  
» for \$n2.lemma give \$1, count(), sort by \$2 desc, \$1
- ▶ a-node \$n1:= [ lemma = "forma", child \$n2:= [ conll/depel = "amod" ] ]  
» for \$n2.lemma give \$1, count(), sort by \$2 desc, \$1

## Lemmatization & Part-of-Speech Tagging

### Textual Resources

Lemmatized Corpora for Latin

### Tools and Hands-on

Tools for lemmatization & POS Tagging

Try it yourself on the provided text(s)!

### Latin in the Semantic Web

The LiLa Knowledge Base

The TextLinker

### Querying Lemmatized Resources

Latin Treebanks in Universal Dependencies

Interlinked Lexical and Textual Resources in LiLa

## SPARQL Access Point

<https://lila-erc.eu/sparql/>

# Thanks!

Get in touch



## LiLa: Linking Latin

Università Cattolica del Sacro Cuore  
CIRCSE Research Centre



[info@lila-erc.eu](mailto:info@lila-erc.eu)



<https://github.com/CIRCSE>



<https://lila-erc.eu>



@ERC\_LiLa



Largo Gemelli 1, 20123 Milan, Italy



This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme - Grant Agreement No. 769994.