



UNIVERSITÀ
CATTOLICA
del Sacro Cuore



The LiLa Knowledge Base in a Nutshell

Eleonora Litta

LiLa Tutorial

Mobility & Humanities: Digital Experiences and Tools (2021)

1 December, Università degli studi di Padova



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.

Introduction and fundamentals

- LiLa: mission and architecture

LiLa now!

- Lemma Bank and Lexical Resources

- Textual Resources

- To sum up

Tutorials

- What, How and Who

- Text Linker

We have built and collected (for Latin and other languages):

- ▶ Textual Resources
- ▶ Lexical Resources
- ▶ NLP Tools

Scattered and unconnected

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

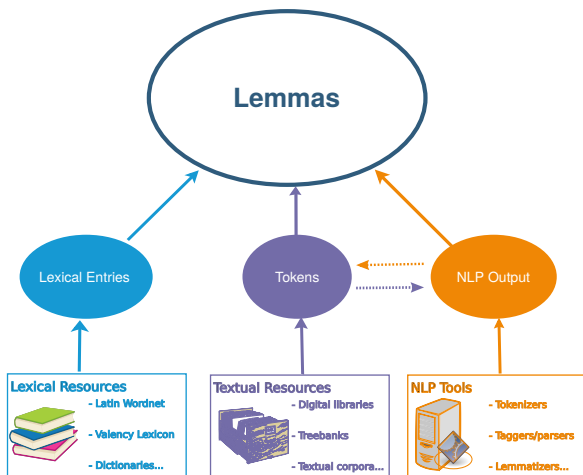


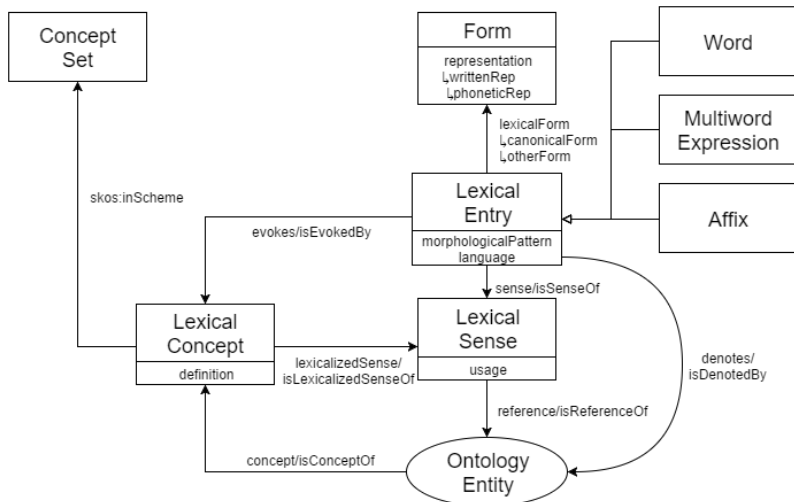
Infrastructure



Interoperability

- ▶ Use URIs for things (e.g. an entry in a lexicon, a token in a corpus) - FINDABLE
- ▶ Use HTTP URIs to allow people (and machines) to look up things - ACCESSIBLE
- ▶ Use web standards to represent/query (meta)data, such as RDF and SPARQL - INTEROPERABLE
- ▶ Include links to other URIs - INTEROPERABLE
- ▶ metadata are released with a clear and accessible data usage license - REUSABLE





LiLa is based on an ontology made of:

- ▶ **Individuals:** instances of objects (one specific token, lemma etc.)
- ▶ **Classes:** types of objects/concepts (token, lemma, PoS etc.)
- ▶ **Data properties:** attributes that objects can/must have (morphological features for lemmas/tokens)
- ▶ **Object properties:** ways in which classes and individuals can be related to one another: RDF triples.

Labels from a restricted vocabulary of knowledge description:

hasLemma, hasPoS

Each component of the ontology is uniquely identified through a URI.

Lemma *admiror* 'to admire, to respect'

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

- ▶ Lemma Bank
- ▶ 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
- ▶ A bilingual Latin-English Dictionary (Lewis & Short)

LiLa: Lexical Resources. All together now!

Lemma Bank, Derivational Morphology, Etymology and Polarity

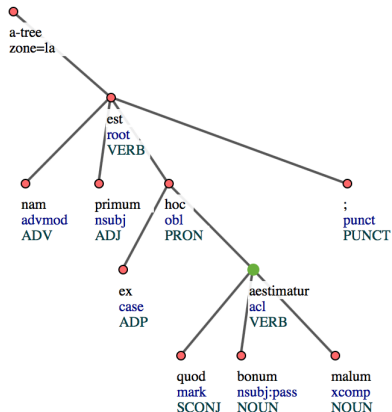


(Annotated) Corpora in LiLa

Source: The *Index Thomisticus* Treebank (UD 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)*

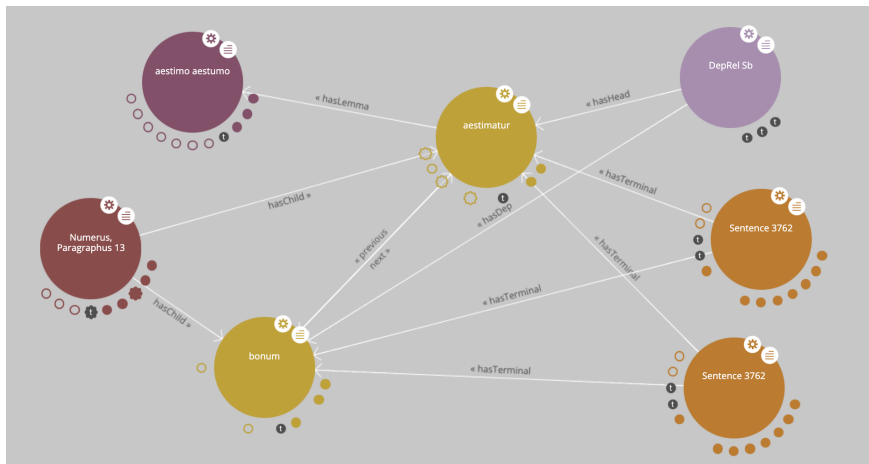


Token *aestimatur*

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

Texts, tokens, relations and lemmas

Phenomena and noumena



► Corpora

- ✓ Index Thomisticus Treebank (*Summa contra Gentiles*): ca. 450,000 nodes
- ✓ Dante Search (700th death anniversary): ca. 46,000 tokens
- ✓ Querolus sive Aulularia: ca. 17,000 tokens
- PROIEL and LLCT treebanks
- Computational Historical Semantics, LASLA and CroALa Corpora

► Lexica

- ✓ Word Formation Latin: ca. 46,000 lemmas (Classical Latin)
- ✓ Etymological dictionary of Latin & the other Italic Langs.: ca. 1,400 entries
- ✓ LatinAffectus: ca. 2,300 entries
- ✓ Index Graecorum Vocabulorum in Linguam Latinam: ca. 1,800 entries
- ✓ Latin WordNet: ca. 1,000 manually checked entries
- ✓ Latin Vallex 2.0: Valency Lexicon
- Lewis & Short Dictionary

► NLP tools

- ✓ LEMLAT (lemma bank): ca. 150,000 lemmas

► TOTAL: approximately 10 million triples

Query Interface, Triplestore

- ▶ Query interface
- ▶ Triplestore

Linguistic Resources. Corpora

- ▶ Index Thomisticus Treebank
- ▶ UDante
- ▶ Querolus sive Aulularia
- ▶ Liber Abbaci

Linguistic Resources. Lexica

- ▶ Word Formation Latin
- ▶ Etymological Dictionary of Latin & the Other Italic Languages
- ▶ LatinAffectus
- ▶ Index Graecorum Vocabulorum in Linguam Latinam
- ▶ Latin WordNet
- ▶ Latin Vallex 2.0

- ▶ **WHAT:** to show the workflow we employ to connect a linguistic resource for Latin to the LiLa Knowledge Base, and to demonstrate the way in which LiLa can be queried.
- ▶ **HOW:** to teach participants how to perform automatic lemmatisation and RDF-isation (i.e. format conversion) in order to link a Latin text to LiLa.
- ▶ **WHO-1:** anyone who wishes to publish Latin texts on the web.
- ▶ **WHO-2:** anyone interested in the different aspects involved in the construction of a Linguistic Linked Open Data knowledge base.

LiLa reflects the annotation granularity of the resources it connects

No data enrichment or further analysis is performed
...but we can help you to enrich your (meta)data

To enter the LiLa Knowledge Base, a textual/lexical resource must be:

- ▶ Lemmatised
- ▶ Part-of-Speech tagged (ideally, using the Universal Dependencies tagset)
- ▶ Online!



Figure: LiLa's Text Linker


LILA: TEXT LINKER (β)

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.

LILA KNOWLEDGE BASE LINKING



exact match
ambiguous match
no match

Click a token to show linked data

Form: basia

Lemma: basium - Upos: NOUN

Data from LemmaBank:

Linked to LiLa [LiLaLemma:91394](#)

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

Copyright © LiLa ERC 2020

Figure: Text processed against the LiLa Knowledge Base

<http://lila-erc.eu:8080/LiLaTextLinker/>

- ▶ LiLa Lemma Bank query interface: <https://lila-erc.eu/query/>
- ▶ LiLa SPARQL interface: <https://lila-erc.eu/sparql/>
- ▶ Repository for the tutorial:
<https://github.com/CIRCSE/Tutorials/tree/main/MobiLab>
- ▶ Collaborative notebook: <https://tinyurl.com/2p83buww>
- ▶ TextLinker (!!Beta Version!!):
<http://lila-erc.eu:8080/LiLaTextLinker/>
- ▶ Repository of SPARQL queries:
<https://github.com/CIRCSE/SPARQL-queries>

LiLa: Linking Latin

Università Cattolica del Sacro Cuore
CIRCSE Research Centre



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.