# Interlinking Corpora Through Lemmas

The LiLa Knowledge Base
of Interoperable Linguistic Resources for Latin

**Marco Passarotti**

CL 2021
13th - 16th July 2021, University of Limerick, Ireland

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)

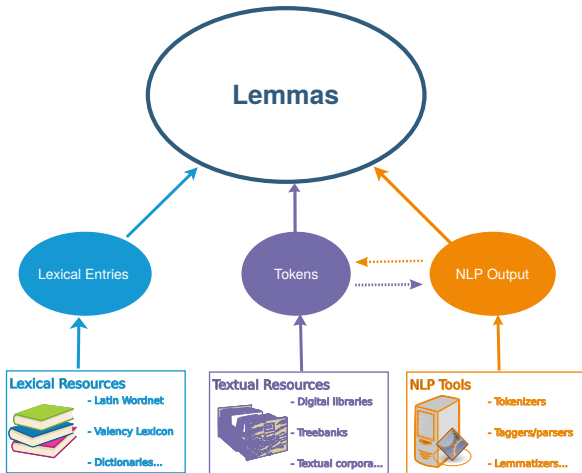## Interlinking as a Form of Interaction



**Infra**structure



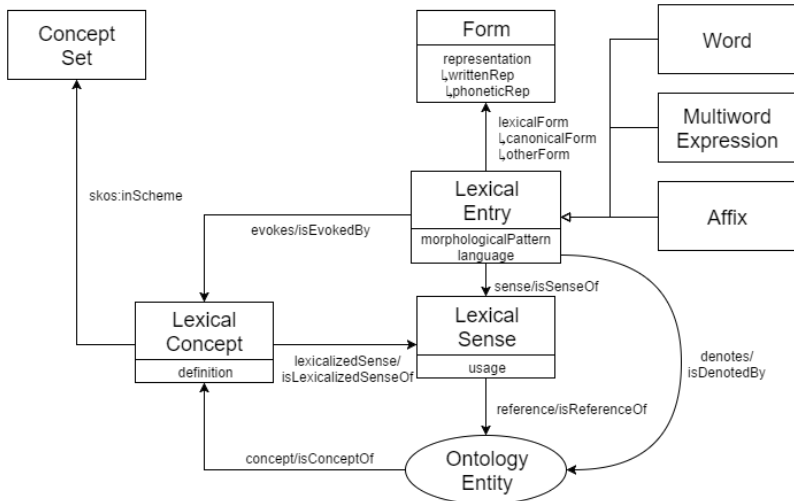**Inter**operability

- ▶ 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

Lemma *admiror* 'to admire, to respect'
`https://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)
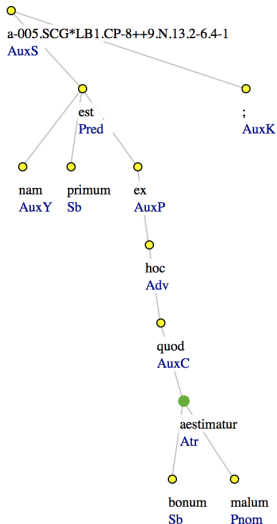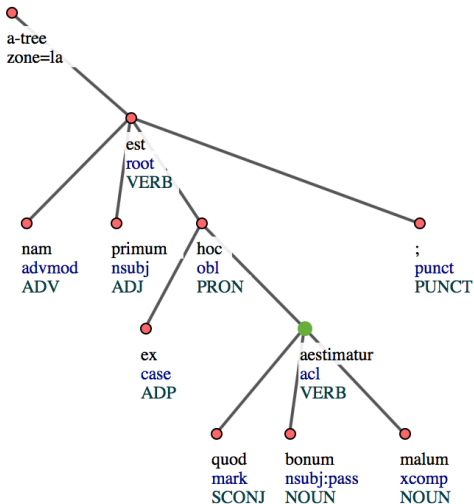► A manually checked subset of the Latin WordNet

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

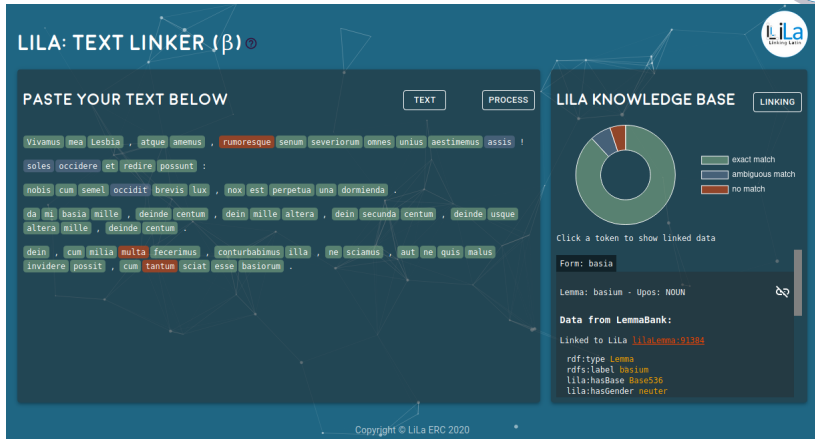*for the first arises because the good **is judged** to be evil;* (Trans. Anton C. Pegis)

Token *aestimatur*

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

# TextLinker
Welcome screen



Figure: LiLa's Text Linker

# TextLinker
## Processed output



Figure: Text processed against the LiLa Knowledge Base

`http://lila-erc.eu:8080/LiLaTextLinker/`

► **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 13 million triples**

### Query Interface, Triplestore and Linker

- ▶ Query interface; Triplestore
- ▶ Linker

### Linguistic Resources. Corpora

- ▶ Index Thomisticus Treebank
- ▶ Dante Search
- ▶ *Querolus sive Aulularia*

### 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

## 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