# Innovation in language resources development
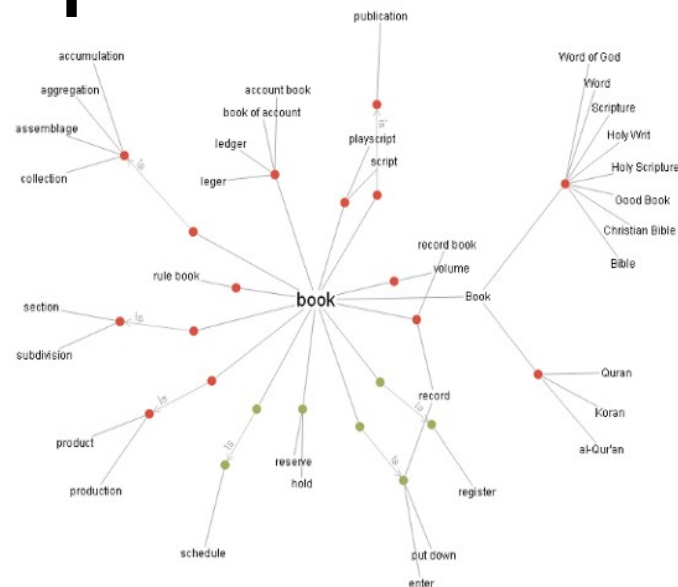
Water flows, rocks remain.
(Romanian proverb)

# Major resources development

Romanian WordNet

– lexical ontology

– 60,000 synsets

– aligned with Princeton WordNet for English and thus with other wordnets aligned to it: e.g. Hindi, Assamese, Bengali, Boro, etc.
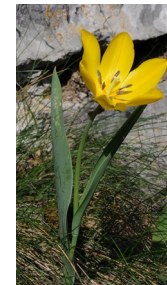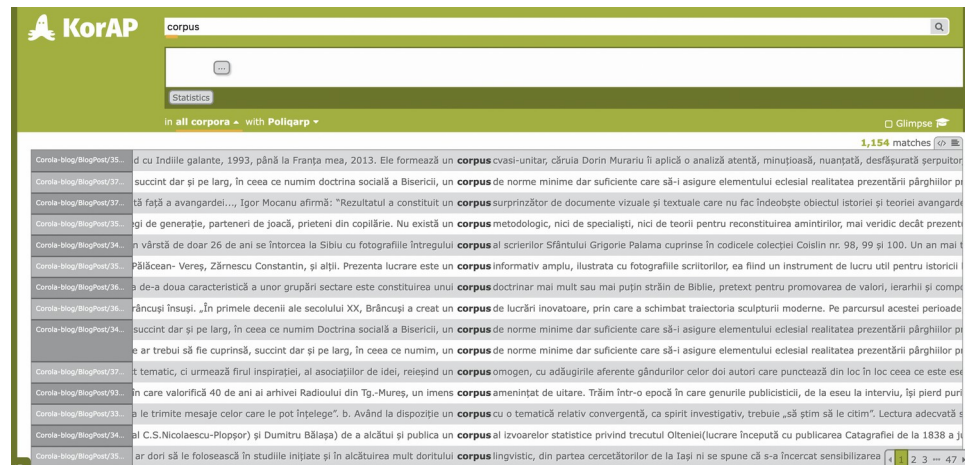
# Corpora: CoRoLa

- Priority project of the Romanian Academy

- Partners:
  - Institute of Computer Science
  - University of Bucharest, Institute for the German Language
  - Texts providers

- 1.2 mill. words

- Bimodal: written, spoken

- diverse wrt styles, domains, subdomains

- several annotation levels

- metadata

- exploitation: KorAP (manual), OCQP (oral component), word embeddings (for downstream applications)

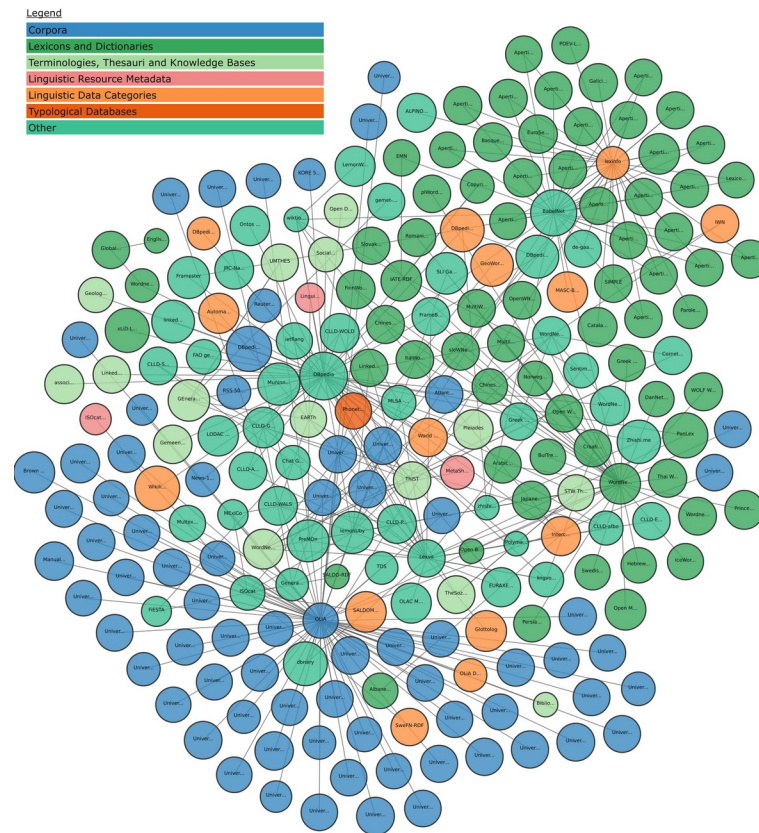**Star of the Carpathians**

# Other corpora



- Oral: ROBINTASC

- NER: LegalNERo, SiMoNERo, MARCELL

- Multiple annotation layers: lexical, morphologic, syntactic, phonetic

- Exploited in various tasks: anonymisation, machine translation, etc.

# Standardization

- Linked Data format

- Nexus Linguarum COST Action

- Data FAIRness:
  – Findable
  – Accessible
  – Interoperable
  – Reusable



Legend
Corpora
Lexicons and Dictionaries
Terminologies, Thesauri and Knowledge Bases
Linguistic Resource Metadata
Linguistic Data Categories
Typological Databases
Other

The Linguistic Linked Open Data Cloud from lod-cloud.net

# Access

- Free

- Sometimes limitations: CoRoLa

- Metadata available in major European Language Technologies hubs:
  - META-SHARE
  - European Language Grid
  - Linked Open Data Cloud

- Data dump available