# The **LexBib** Project

A Corpus and Bibliography of Metalexicography

*David Lindemann*

*Universität Hildesheim*

*UPV/EHU University of the Basque Country*

# LexBib Project

- Goals

- Features, Methods, Workflow

- Partners

- Working plan, Outlook



- Detailed intro: Lindemann, Kliche & Heid 2018 (Euralex Ljubljana)

# State of the Art and Goals

**Situation now**

- Bibliographies: Unlinked islands
  - Publishers / topics

- Noisy metadata sets
  - Automatically harvested

- Language barrier
  - Mutual invisibility

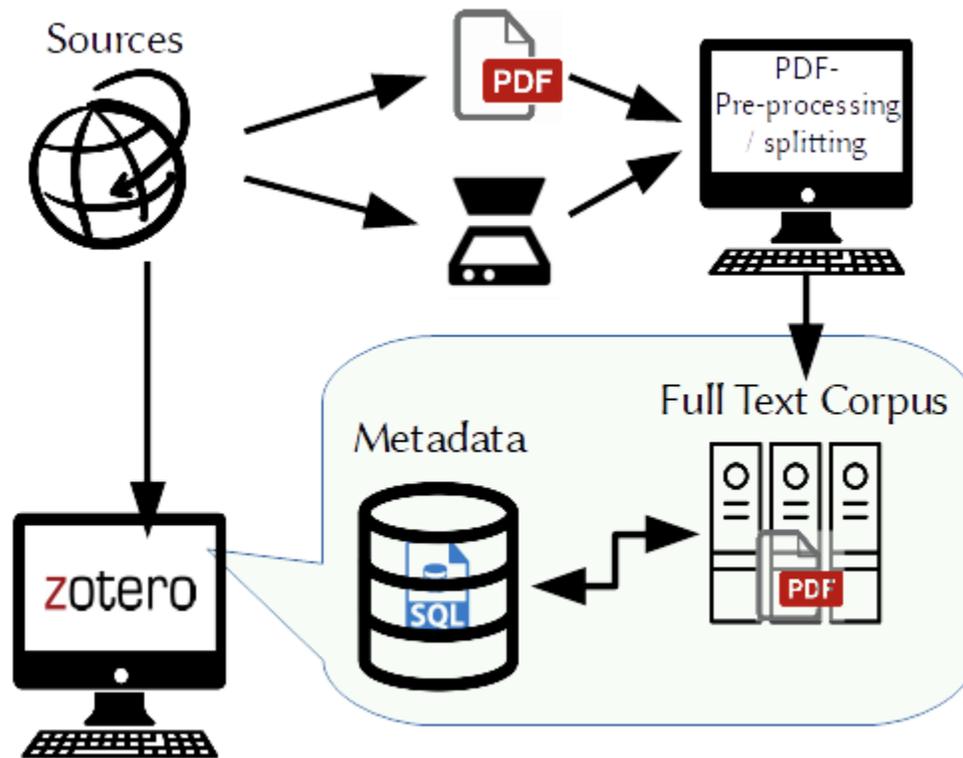- Unconsistent indexation
  - Keywords / subject headings

**LexBib Goals**

- Unified territory
  - But: only relevant content

- Complete and correct metadata
  - manually validated

- Multilingual access structure
  - Multilingual Domain Ontology
  - Additional content-describing metadata
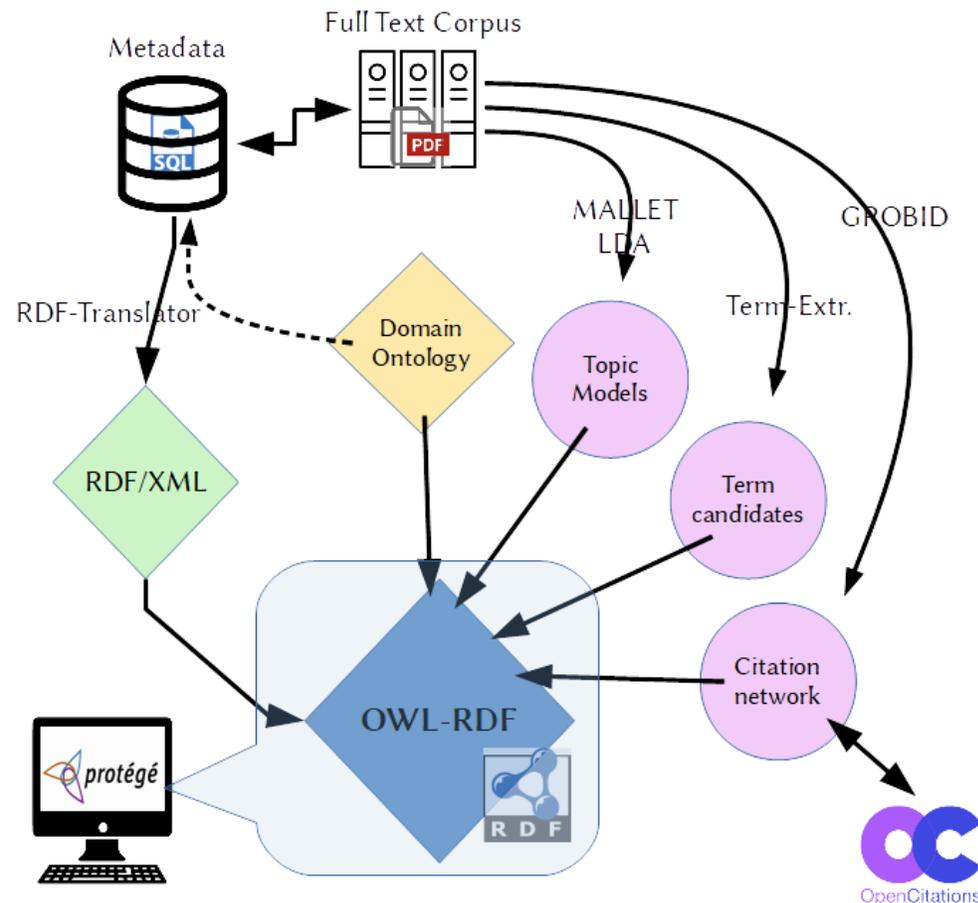
# LexBib online reference platform

- Advanced search & browsing functions

- Access through domain ontology nodes

- Access through citation network
  - Visualisation of network and clusters

- Author profiles
  - Author name disambiguation, additional info

- Link to full text
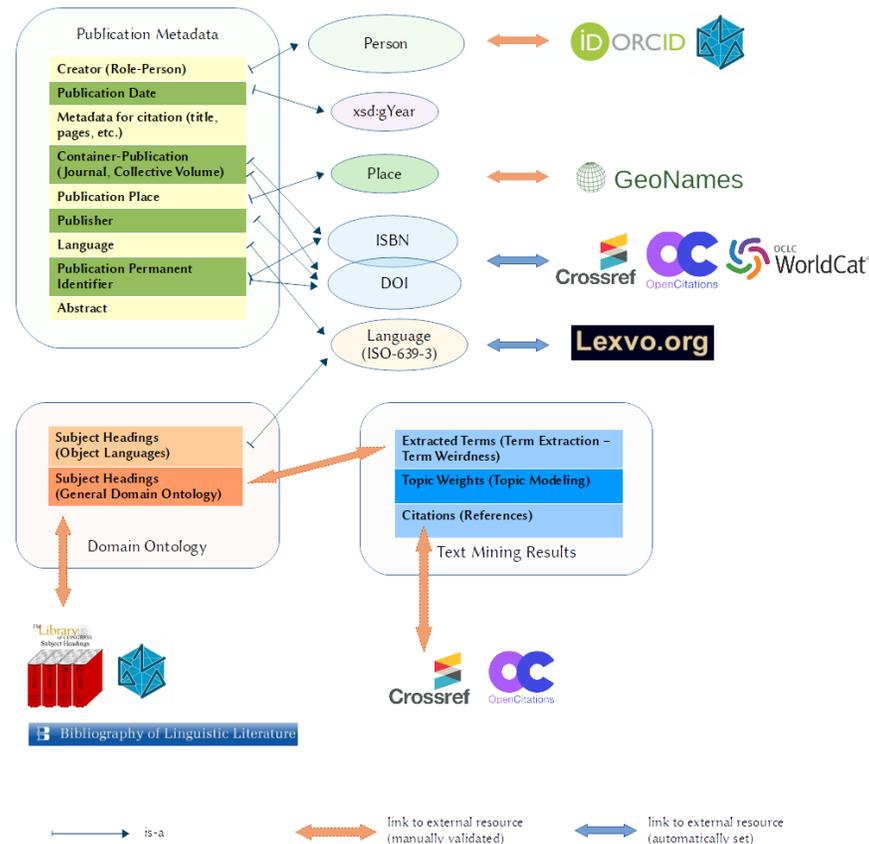  - …on publishers' etc. platforms

# Working Pack 1

- **Collection and validation of Metadata**
  - If necessary
    - completion/correction
- **Collection of full texts**
  - If necessary
    - Scan
    - OCR
    - Split (collective volumes)

# Working Pack 2



- **Extraction from full texts**
  - **Term / keyword candidates**
    - Weirdness / reference corpora
    - Tools: Rösiger, Heid et al. (IMS Stuttgart)
  - **Topic Models**
    - Topic Weights (for recommender)
    - Tools: McCallum: MALLET
  - **Citation relations**
    - Citation network, clusters
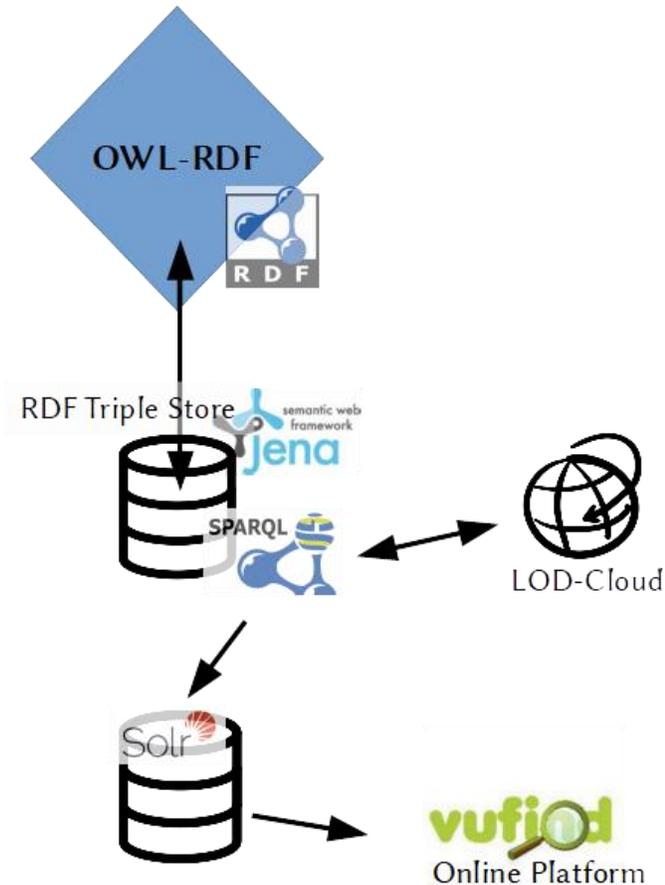    - Tools: Romary, Lopez et al. (GROBID)

# Working Pack 3



- Domain Ontology
  - Controlled Vocabulary of subject headings
    - Root node "Lexicography", linked to LCSH, DNB, Wikidata...
  - Multilingual labels
    - EN, DE, ES, ...
  - Sources: Existing classifications
- Linked Open LexBib
  - Use of existing element sets
  - Links to existing resources

# Working Pack 4



- **Open Zotero Group**
  - zotero.org/groups/LexBib
- **Linked Open LexBib**
  - SPARQL Endpoint
- **LexBib Online Platform**
  - Implementation with *vufind*

# Working Pack 5

- Dissemination, Community
  - Euralex, eLex conferences
    - Feedback
    - Call for contributions
    - Curation of author pages

- Process Metadata
  - Predictions for similar projects on a larger scale:
  - *"How much manual work is necessary for a comparatively noise-free and complete resource?"*

# Partners

**DFG-Proposal (in prep.): Partners**

- Uni Hildesheim: Lead, WP 1-5
- BBAW: WP 2, 5
- IDS: WP 1, 4, 5

**Financial and other support**

- Uni Hildesheim: Preliminaries
- Elexis: Preliminary work WP 1, 3
- Euralex
- EMLex
- Lexicografía UDC
- LOC-DB / VZG
- De Gruyter

# Working Plan

- 2019 [Uni Hildesheim, Elexis]
  - Metadata: Euralex-Dykstra, Obelex-Meta, Wiegand, Córdoba, Ahumada
  - Full Texts: English 2000-2018
  - Domain Ontology: First proposal

- 2020-22 [DFG?]
  - Metadata: All European languages, 1971-2018
  - Full texts: German, Spanish 2000-2018 (and more?)
  - WP2-5

- In a near future
  - Other languages, older full texts

# Thank you!
# Questions?

david.lindemann@uni-hildesheim.de

http://zotero.org/groups/LexBib