

Open Linguistics Working Group (OWLG)

Christian Chiarcos
chiarcos@uni-frankfurt.de



Open Knowledge Foundation (OKFN, <http://okfn.org>)



- non-profit organization
- founded in 2004
- promote open knowledge in all its forms
 - e.g., publication of government data (UK, US)
- provide infrastructural support for several working groups



OKFN Open Linguistics Working Group (OWLG)

- founded in Oct 2010 in Berlin, Germany
- open network of individuals interested in
 - linguistic resources and/or
 - their publication under open licenses
- multi-disciplinary
 - NLP/CL, typology/language documentation, SW, ...
- infrastructure
 - mailing list, web site/blog, wiki
 - <http://linguistics.okfn.org>

OWLG goals

(<http://linguistics.okfn.org>)

1. **Promote open data** in relation to language data
2. Point of **reference and support** for open linguistic data
3. **Facilitate communication** between researchers that use, distribute, or maintain open linguistic data
4. **Mediate between providers and users** of technical infrastructures
5. Build and maintain an **index of open linguistic data sources**
6. Assemble **best-practice guidelines and use cases** concerning creating, using and distributing data
7. Gather **information on legal issues**

OWLG goals

(<http://linguistics.okfn.org>)

1. **Promote open data** in relation to language data
2. Point of **reference and support** for open linguistic data
3. **Facilitate communication** between researchers that use, distribute, or maintain open linguistic data
4. **Mediate between providers and users** of technical infrastructures
5. Build and maintain an **index of open linguistic data sources**
6. Assemble **best-practice guidelines and use cases** concerning creating, using and distributing open linguistic data
7. Gather **information on** open linguistic data

these aspects are
specifically well developed

OWLG activities

- mostly point-to-point cooperations between individual members
- regular telcos/meetings
- workshops -> building an interdisciplinary community
 - colocated with larger events of different communities
 - Linguistics Track of the OKCon, June 2011, Berlin, Germany
 - Linked Data in Linguistics
 - March 2012, Frankfurt/M., Germany -> linguistics / NLP
 - Sep 2013, Pisa, Italy -> academic linguistics
 - May 2014, Reykjavik, Iceland -> NLP/semantics
 - MLODE-2012, Sep 2012, Leipzig, Germany -> IT
 - Linked Data in Linguistic Typology, Sep 2013, Leipzig, Germany



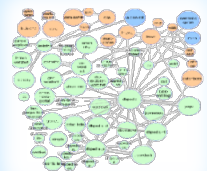
OWLG activities

- point-to-point cooperations between individual members
- regular telcos/meetings
- workshops -> building an interdisciplinary community
 - keeping ties with other communities & projects
 - e.g., Cyberling, W3C OntoLex, ACL SIGANN/SIGLEX
 - e.g., MPI-EVA, LOD2, LIDER, QTLeap
- joint publications and presentations
- building and maintaining the Linguistic Linked Open Data (LLOD) [sub-]cloud

LLOD cloud

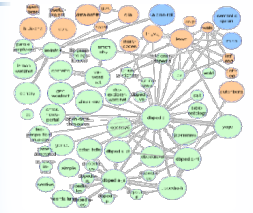
- a collection of linguistic resources
 - ❑ published under open licenses
 - ❑ as linked data
 - ❑ decentralized developed and maintained
 - ❑ meta data at <http://datahub.io>
 - => cloud diagram
 - ❑ developed as a community effort in the context of the Open Linguistics Working Group of the Open Knowledge Foundation

next:
LLOD 2011-2014



Building the Cloud: 2011

A sketch from a table napkin



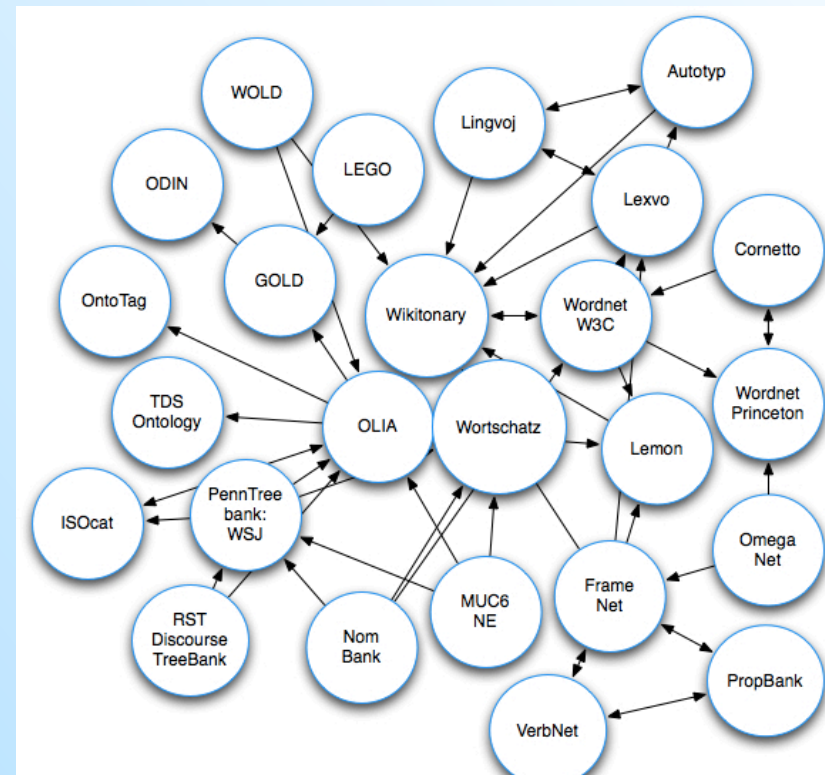
- initially, we maintained a list of open **or** representative resources

- in Jan 2011, we marked possible synergies

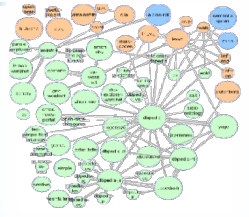
- merely a vision

- includes non-open resources as placeholders for other resources to come
 - not physically realized

- a strong metaphor brought to a new community

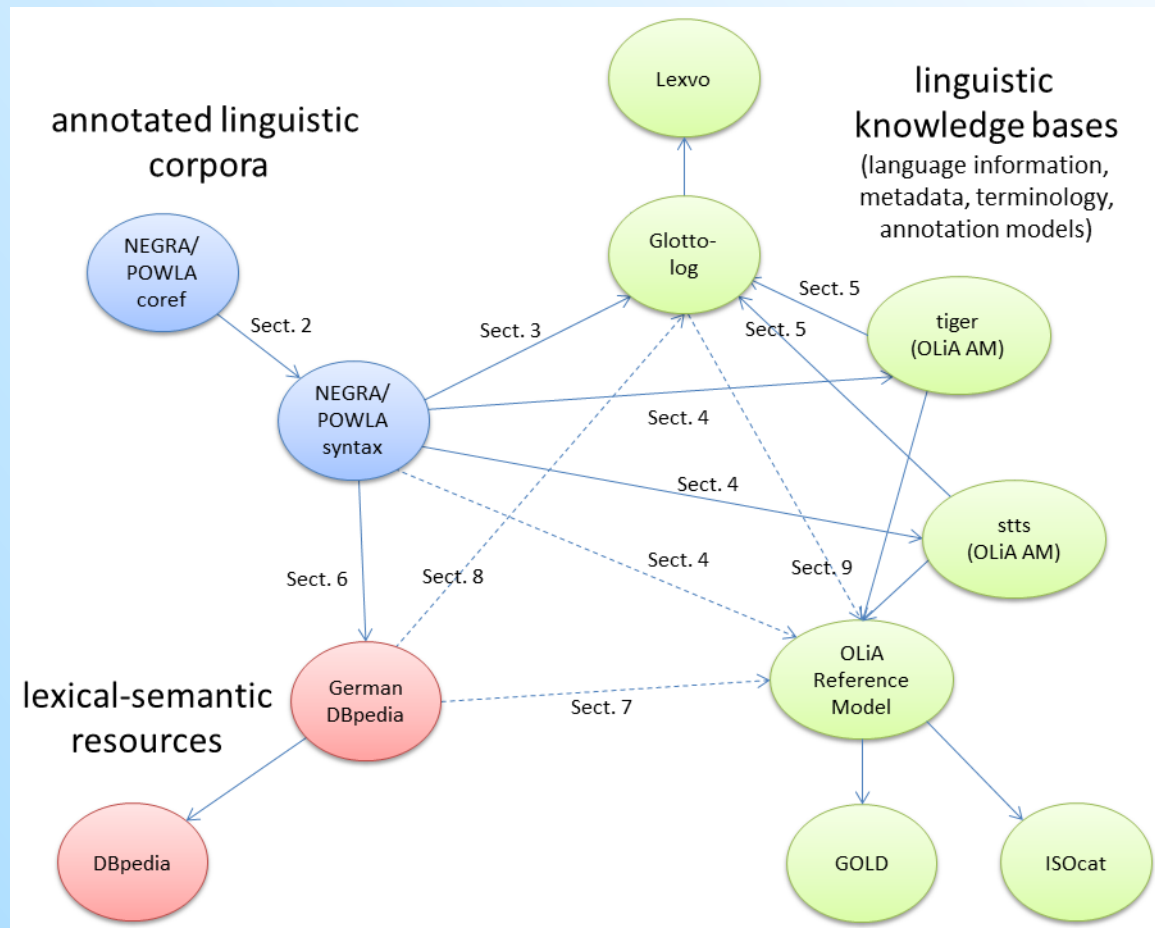
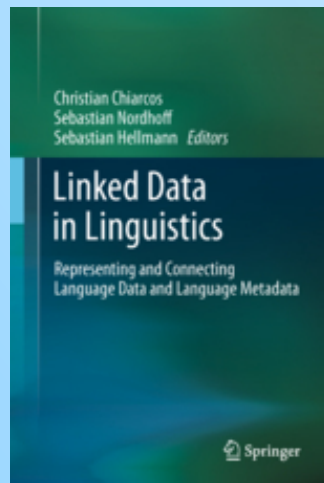


Chiarcos, Hellmann & Nordhoff „Linking Linguistic Resources“ (2012)

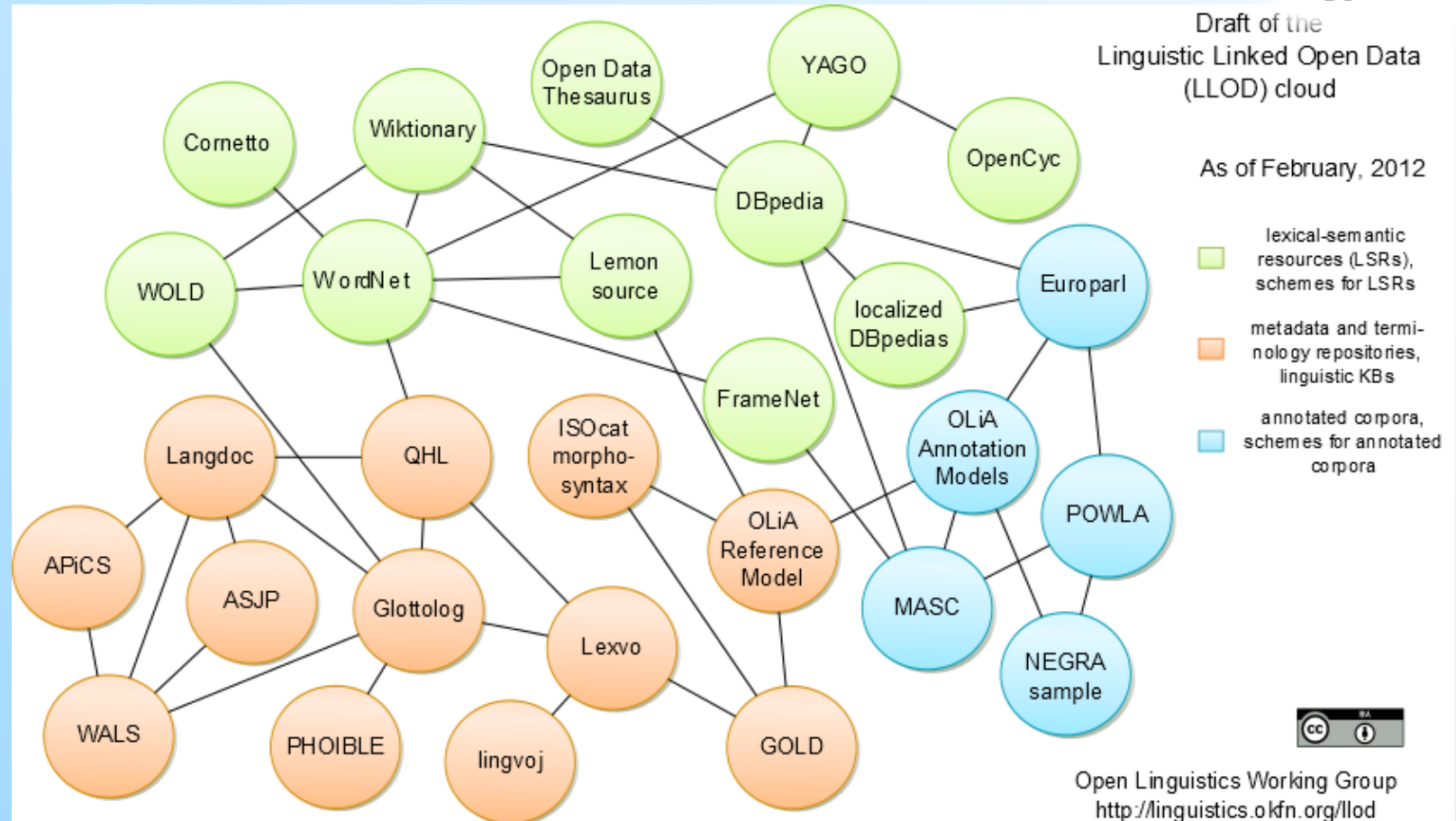
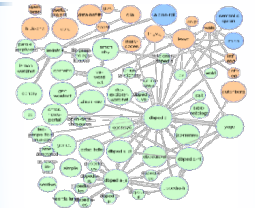


Closing chapter of the LDL-2012 companion volume

- **hypothetical** linking for selected data sets from NLP, SW and typology described in the book



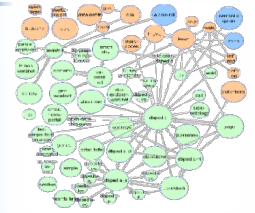
Drafting the Cloud: LREC-2012



„draft status“

hand-crafted, including resources whose RDF conversion and linking was **suggested**, not yet performed at the time

Building the Cloud: MLODE-2012

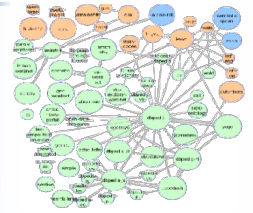


- Multilingual Linked Open Data for Enterprises
 - ❑ goal: build the first instance of the LLOD cloud
 - ❑ workshop & hackathon
 - authors were encouraged to provide data
 - data conversion, metadata update at <http://datahub.io>
- automatically generated diagram
 - ❑ Richard Cyganiac's converter scripts



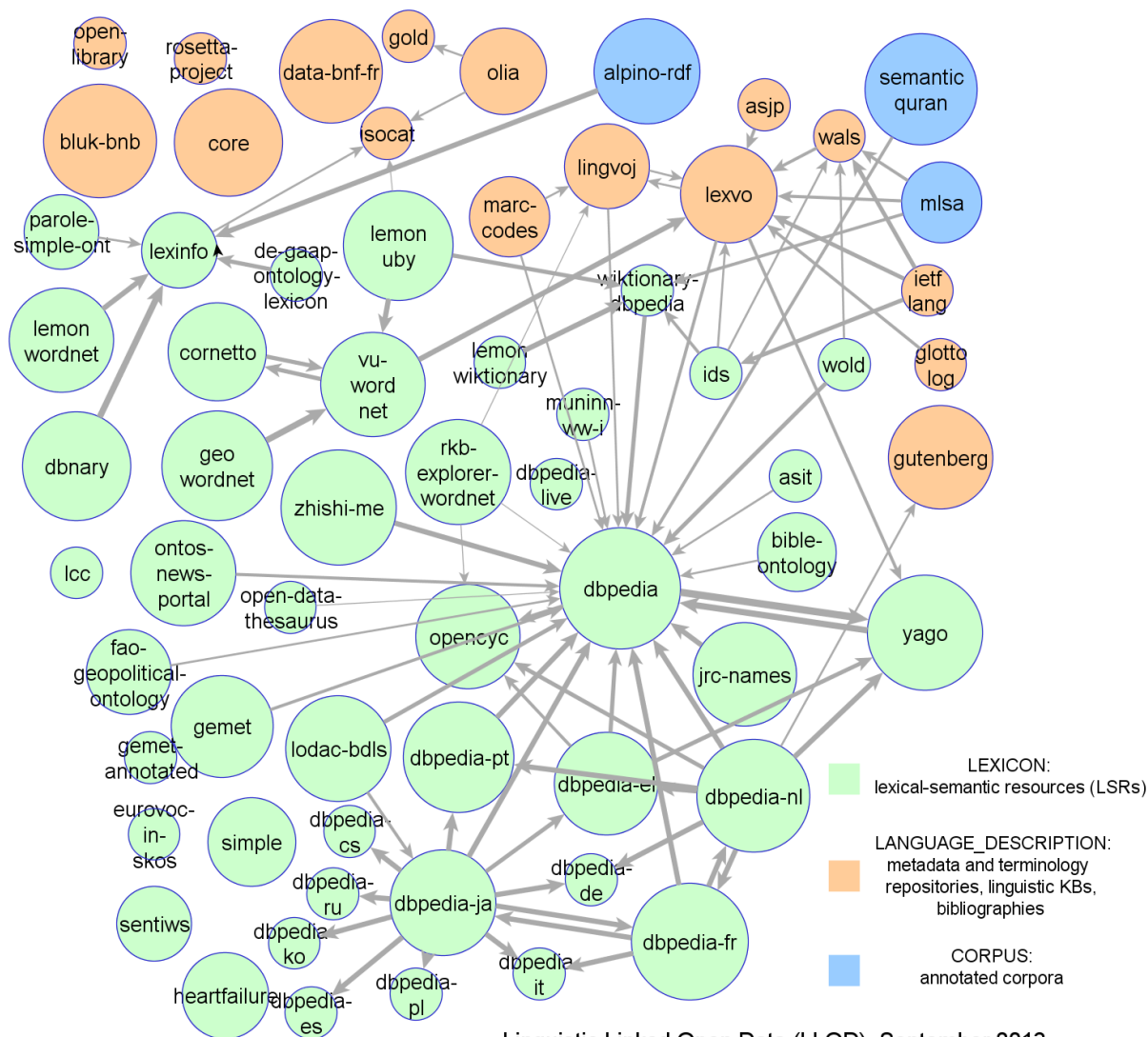
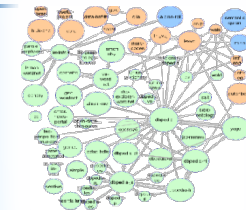


Building the Cloud: 2013+



- MLODE data post-proceedings
 - Special issue of the Semantic Web Journal
 - Preparation of additional data sets in the process
 - e.g., lemonUby (Eckle-Kohler et al., accepted)
- Linked Data in Linguistic Typology, Aug 2013
 - additional potential datasets
 - lexical databases of Austronesian languages
 - a database of syllable structures
- Intensified community work

Building the Cloud: Sep 2013



Linguistic Linked Open Data (LLOD), September 2013
CC-BY Open Linguistics Working Group (<http://linguistics.okfn.org/llod>)

- more data sets
- not fully linked, yet
- new drawing script
 - by John McCrae & Christian Chiarcos
- manually categorized and colored
 - GraphML

Build Cloud

The graph illustrates a dense network of semantic resources. Central nodes include 'dbpedia' and 'yago'. Other prominent nodes are 'lexinfo', 'gutenberg', 'semantic quran', 'lingvoj', 'lexvo', 'cornetto', 'greek-wordnet', 'geowordnet', 'EARTH', 'agrovo-skos', 'pleiades', 'lodac-balls', 'dbpedia-es', 'dbpedia-ko', 'dbpedia-cs', 'dbpedia-it', 'dbpedia-fr', 'dbpedia-de', 'dbpedia-el', 'dbpedia-nl', 'dbpedia-ja', 'dbpedia-pt', 'dbpedia-uk', 'dbpedia-ru', 'dbpedia-sr', 'dbpedia-sl', 'dbpedia-sv', 'dbpedia-tr', 'dbpedia-zh', 'dbpedia-he', 'dbpedia-hi', 'dbpedia-bn', 'dbpedia-ta', 'dbpedia-te', 'dbpedia-mr', 'dbpedia-mn', 'dbpedia-my', 'dbpedia-th', 'dbpedia-vi', 'dbpedia-id', 'dbpedia-fa', 'dbpedia-ar'. The nodes are color-coded: blue for 'Build', green for 'Cloud', and orange for 'Data'.

- | | |
|---|--|
| <p>LEXICON:
lexical-semantic resources (LSRs)</p> <ul style="list-style-type: none"> ○ general semantic knowledge bases ● lexical resources with grammar information | <p>METADATA:
information about language and language resources</p> <ul style="list-style-type: none"> ○ information about language resources (incl. bibliography) ● linguistic terminology repositories |
| <p>CORPUS:
collections of language samples</p> <ul style="list-style-type: none"> ○ annotated corpora | <p>● databases of language features (e.g., from typology)</p> |

created in preparation of the 3rd Linked
Data in Linguistics Workshop (LDL-2014)

Recent developments

- finalizing LLOD diagram revision
 - for LDL-2014, May 27th, 2014
- harmonizing linguistic resource categories
 - synchronization with MetaShare categories
- adding new resources
 - relevant LREC „Share your resources“ datasets ?
- subsequently enforce further constraints on LLOD „bubbles“
 - open licenses (currently: accessible ~ LOD diagram)
 - well-formedness / meta data check