

david.lindemann@uni-hildesheim.de
fritz.kliche@uni-hildesheim.de
heid@uni-hildesheim.de

LexBib

A corpus and bibliography
of metalexicography



David Lindemann
Fritz Kliche
Ulrich Heid

LexBib Goals

▼ Infrastructural aspect

- ▼ Bibliography of Lexicography and Dictionary Research (Metalexicography)
- ▼ complete, structured, validated
 - ▼ Full text collection (“e-Science Corpus”) > Computational text analytics
 - ▼ Publication metadata collection > Online-Bibliography

▼ Research Aspect

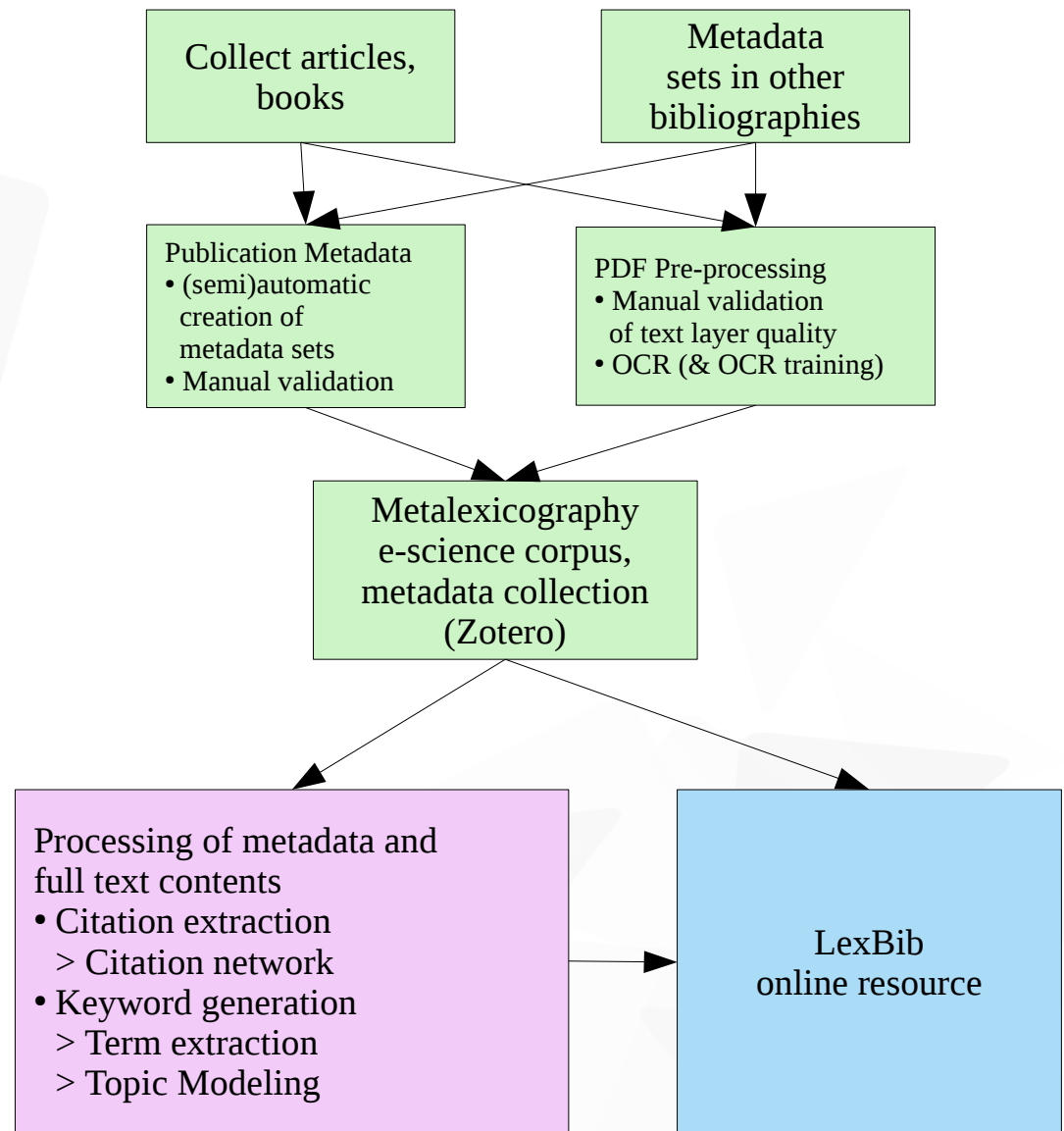
- ▼ Computational methods, applied to full texts
 - ▼ Bibliometrics: Citation Extraction > Citation Network
 - ▼ Topic Modeling, Term Extraction > Keyword Indexation
 - ▼ Tests and Evaluation
- ▼ Results: Additional publication metadata, to be included in online bibliography

▼ Motivation

- ▼ Existing bibliographies far from state-of-the-art resources available for other disciplines
- ▼ Hildesheim and IDS: Existing manually curated metadata collections and full text corpora as starting point
- ▼ Relatively small discipline > Feasible for hand-validation and evaluation (predictions for the application of similar workflows in broader domains)

LexBib project: Main Modules

- ▼ **Green:** creation of e-science corpus, metadata and PDF pre-processing
 - ▼ Current LexBib collection: 2,500 items (July 2018)
- ▼ **Red:** Computational text analysis:
 - ▼ Creation of additional metadata
 - ▼ Citation Network
 - ▼ Keyword generation
- ▼ **Blue:** Creation of online bibliography
 - ▼ Models followed: DBLP, ACM



LexBib project: Module 1

▼ **Green:** creation of e-science corpus, metadata and PDF pre-processing

▼ 1st iteration: Spring 2018

▼ 2,057 items

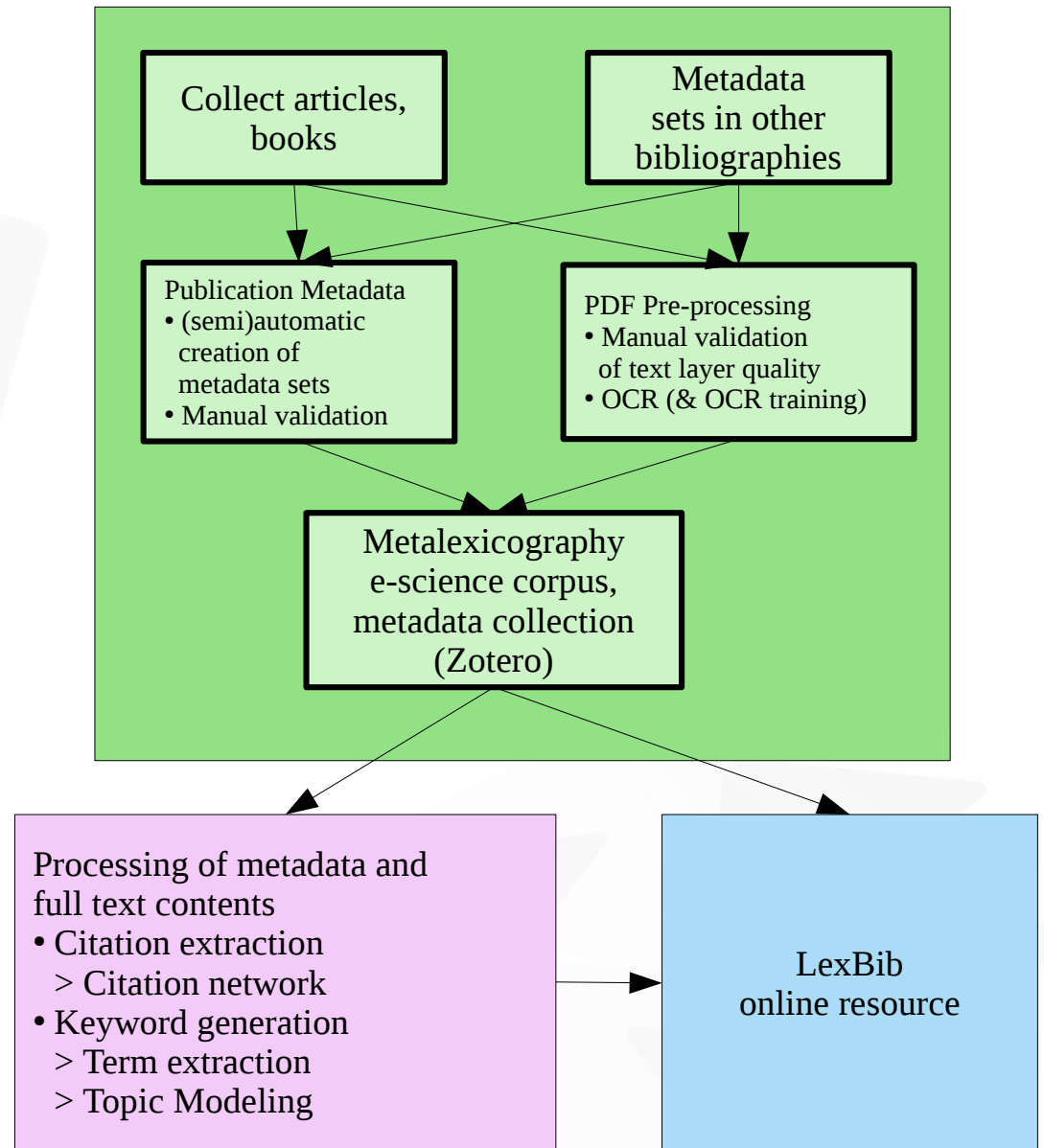
▼ English texts 2000-2017

▼ Core item types

▼ Journal articles

▼ Proceedings (eLex, Euralex)

▼ Handbook articles

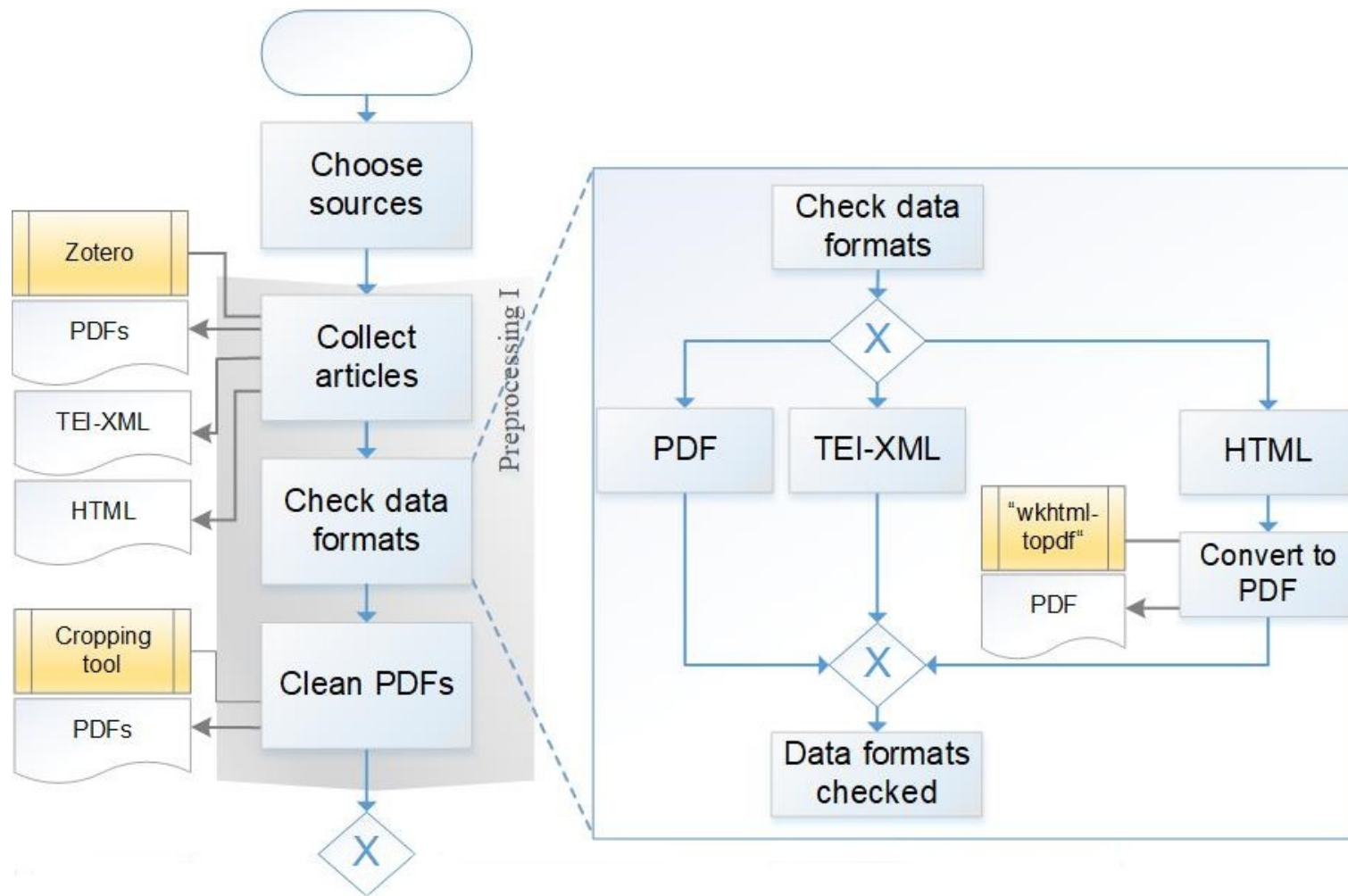


e-science corpus building: metadata & full text collection

- Collection of full texts and metadata using “Zotero”
 - From some online repositories (bulk) download of metadata and full texts possible
 - Zotero cloud: Enables collaborative editing of collections
 - Manual metadata validation/editing in Zotero GUI
 - Manual assessment of PDF quality
 - PDF manipulation: cropping of non-relevant content; running titles, page numbers, boilerplate
 - Re-digitization: OCR training

The screenshot displays the Zotero interface. On the left, a sidebar shows a hierarchical view of collections, including 'COLLECTIVE_VOLS', 'GRANGER_PAQUOT_2012', 'HANDBOOK_HSK_5_4', 'HANDBOOK_ROUTLEDGE', 'CONFERENCES', 'CONF_ELEX', 'CONF_EURALEX', 'CONF_GLOBALEX', 'CONF_LD', 'JOURNALS', 'JOURNAL_ASIALEX', 'JOURNAL_DSNA', 'JOURNAL_IJL', 'JOURNAL_LEXICOGRAPHI...', and 'JOURNAL_LEXIKOS'. The main pane shows a list of items with columns for title, author, and year. The item 'Access to Multiple Lexical Resources at a Stroke: Integrating Dictionary, Corpus and Wordnet Data' by Trap-Jensen (2010) is selected. The right pane shows the detailed metadata for this item, including 'Eintragsart: Konferenz-Paper', 'Titel: Access to Multiple Lexical Resources at a Stroke: Integrating Dictionary, Corpus and Wordnet Data', 'Autor: Trap-Jensen, Lars', 'Herausgeber: Granger, Sylviane' and 'Paquot, Magali', 'Datum: 2010', 'Titel des Konferenzbandes: Proceedings of eLex 2009, Louvain-la-Neuve, 22-24 October 2009', 'Name der Konferenz: eLexicography in the 21st Century: New Challenges, New applications', 'Ort: Louvain-la-Neuve', 'Verlag: UCL Presses', and 'Band'.

Processing Pipeline: Data collection



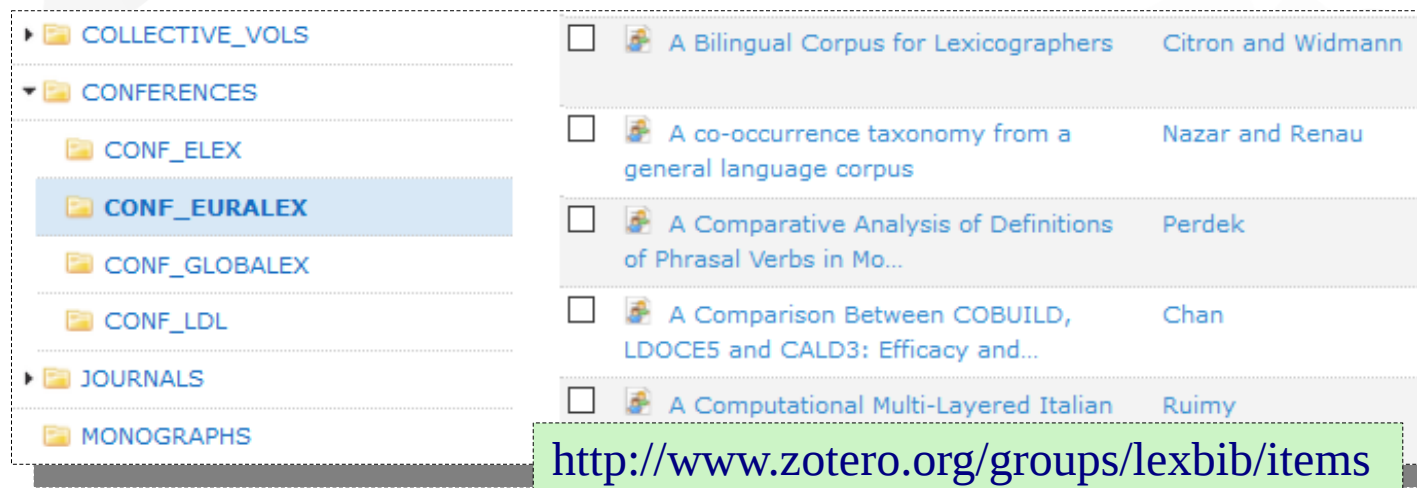
LexBib publication metadata

▼ Hand validated

- ▼ Standard metadata (for citations), according to item type
- ▼ Unique Identifiers: DOI, ISBN, ORCID
- ▼ Item Relations “is_review_of“, “is_reviewed_by“
- ▼ Item Relations “is_container_of“, “has_container“
- ▼ **Red module**: Item Relations “citing“ and “is_cited_by“

▼ Compare and merge

- ▼ LexBib_v1 (Lindemann et al. 2018)
- ▼ IDS Obelex-Meta (Möhre 2016)
- ▼ Euralex-Bib (Dykstra & Hartmann)
- ▼ WLWF-Bib (Kammerer, Gouws et al. in prep.)
- ▼ Wiegand’s „Internationale Bibliographie“ (2006-2014)
- ▼ Córdoba Rodríguez (2003)



The screenshot shows the Zotero interface for the LexBib group. On the left, a sidebar lists item types: COLLECTIVE_VOLS, CONFERENCES, CONF_ELEX, CONF_EURALEX (highlighted), CONF_GLOBALEX, CONF_LD, JOURNALS, and MONOGRAPHS. The main area displays a list of items with checkboxes, titles, and authors:

Item Type	Title	Author(s)
<input type="checkbox"/>	A Bilingual Corpus for Lexicographers	Citron and Widmann
<input type="checkbox"/>	A co-occurrence taxonomy from a general language corpus	Nazar and Renau
<input type="checkbox"/>	A Comparative Analysis of Definitions of Phrasal Verbs in Mo...	Perdek
<input type="checkbox"/>	A Comparison Between COBUILD, LDOCE5 and CALD3: Efficacy and...	Chan
<input type="checkbox"/>	A Computational Multi-Layered Italian	Ruimy

At the bottom of the screenshot, the URL <http://www.zotero.org/groups/lexbib/items> is displayed in a green box.

Some existing bibliographies

<i>Title</i>	<i>Scope (years)</i>	<i>Scope (domains)</i>	<i>Scope (languages)</i>	<i># Items</i>	<i>Format</i>
LexBib Testset	2000-2017	Metalex	English	2,056	Structured database
EURALEX Bibliography	1600-2010	Lex/Metalex	Multiple	1,325	Unstructured (pub. as Wiki)
Obelex-meta	1982-2017	Metalex	Multiple	ca. 2,000	Structured database
WLWF	1420-2016	Lex/Metalex	Multiple	2,370	Unstructured (pub. as PDF)
Wiegand	1850-2014	Lex/Metalex	Multiple	33,339	Unstructured (pub. as PDF)
Hartmann	1930-2007	Metalex	Multiple	ca. 570	Unstructured (pub. as PDF)
Córdoba Rodríguez	1940-2003	Metalex	Multiple	10,192	Structured database
Ahumada	1535-2010	Metalex	Mainly Spanish	6,560	Structured database (in progress)

Wiegand: ‘Internationale Bibliographie der Lexikographie‘

- ▼ 33,000 items, 1850 to 2014
 - ▼ All standard item types + press articles; some with comments. Includes ref. to Zgusta 1988
- ▼ Lindemann & Khemakhem 2018 (in prep.): OCR, Parsing, repr. as XML

21923 Voillat, François: Le “Glossaire des patois de la Suisse romande” (GPSR). In: Actes du XVIII^e Congrès International de Linguistique et de Philologie Romanes [...] Tome VII, 1989[↑], 338–345.

15475 Nyhlén, Lars-Olaf: Wie sagt man in Österreich? In: Moderna Språk 68. 1974, 275–281. [Bes. zu: Jakob Ebner: Duden. Wie sagt man in Österreich? Wörterbuch der österreichischen Besonderheiten. Mannheim 1969 (Duden Taschenbücher 8)].

15397

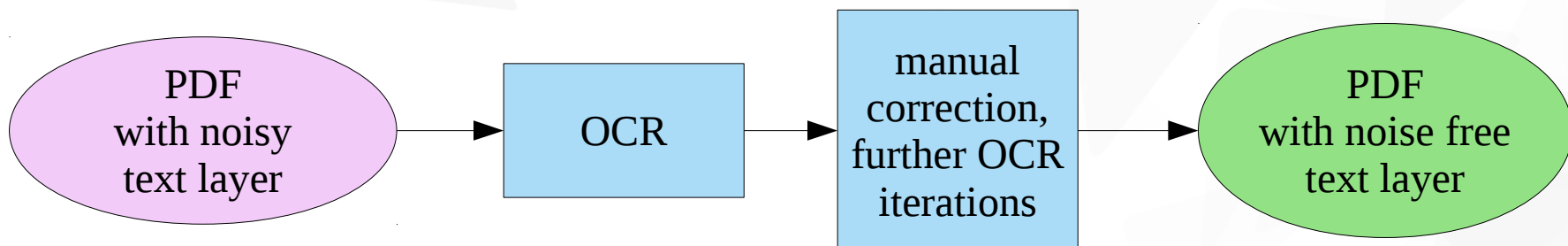
Norden 11.–14. maj 1993. Red. af Anna Garde og Pia Jarvad. Oslo 1993 (Skrifter udgivet av Nordisk Forening for Leksikografi. Skrift nr. 2). [Daraus: 363, 1003, 2755, 3101, 3870, 5222, 5251, 5390, 6058, 6758, 9206, 9743, 10158, 11521, 11553, 11843, 12600, 12940, 12974, 13146, 13399, 14192, 15316, 15345, 15402, 15474, 15651, 15944, 17322, 20818, 21371].

15427 Novikov, L. A.: K probleme omonimii. [Zum Problem der Homonymie]. In: Leksikografičeskij sbornik 1960, 93–102.

Wiegand: ‘Internationale Bibliographie der Lexikographie’

21923 Voillat, François: Le “Glossaire des patois de la Suisse romande” (GPSR). In: Actes du XVIII^e Congrès International de Linguistique et de Philologie Romanes [...] Tome VII, 1989[↑], 338–345.

Original Text Layer	21923.....Voillat,..F.r.a.n.c.i.s.....Le..."Glossaire...des...patois...de...la...Suisse...r.o.m.a.n.d.e." "(GPSR)...In: Actes...du...X.v.ï.ï.Γ...Congres...International...de...Linguistique...et...de...Philologie...R.o.m. .a.n.e.s...[...]. T.o.m.e.VII,..1.9.8.9. ,..3.3.8...3.4.5.
Transkribus OCR	21923.Voillat,·François:·Le·“Glossaire·des·patois·de·la·Suisse·romande”·(GPSR)·In: Actes·du·XVIIIe·Congrès·International·de·Linguistique·et·de·Philologie·Romanes·[...] Tome·VII,·1989 [↑] ,·338-345.
Manual correction	21923·Voillat,·François:·Le·“Glossaire·des·patois·de·la·Suisse·romande”·(GPSR)·In: Actes·du·XVIII ^e ·Congrès·International·de·Linguistique·et·de·Philologie·Romanes·[...] Tome·VII,·1989 [↑] ,·338-345.



Wiegand: Internationale Bibliographie der Lexikographie

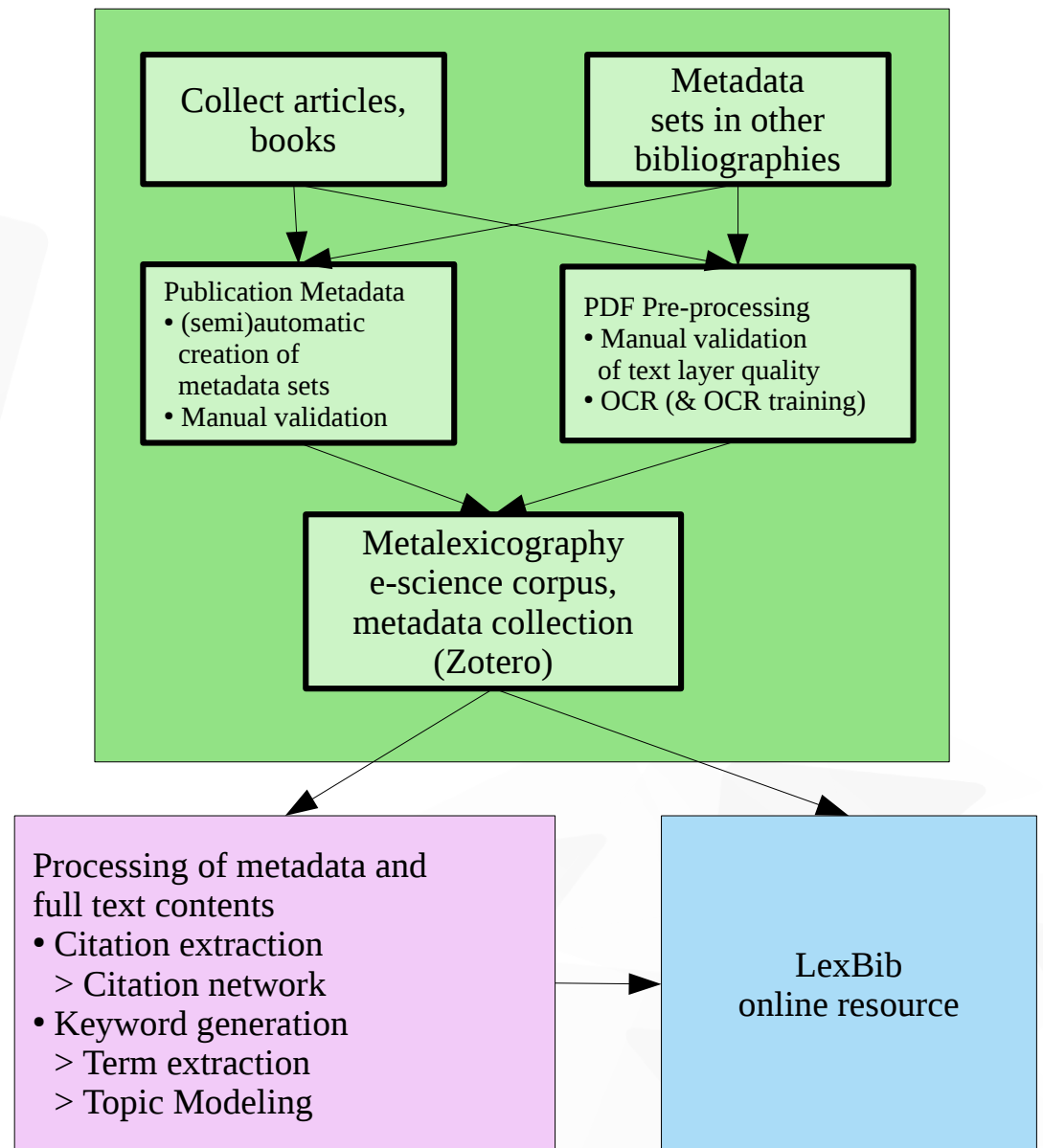
21923 Voillat, François: Le “Glossaire des patois de la Suisse romande” (GPSR). In: Actes du XVIII^e Congrès International de Linguistique et de Philologie Romanes [...] Tome VII, 1989[↑], 338–345.

```
<biblStruct xml:id="Wiegand21923">
  <analytic>
    <author>
      <persName>
        <surname>Voillat</surname><forename>François</forename>
      </persName>
    </author>
    <title level="a">Le "Glossaire des patois de la Suisse romande" (GPSR)</title>
  </analytic>
  <monogr>
    <title level="m">Actes du XVIIIe Congrès International de Linguistique et de
    Philologie Romanes [...] Tome VII</title>
    <imprint><date when="1989"/></imprint>
    <biblScope unit="pp" from="338" to="345">338-345</biblScope>
  </monogr>
</biblStruct>
```

LexBib project: Module 1

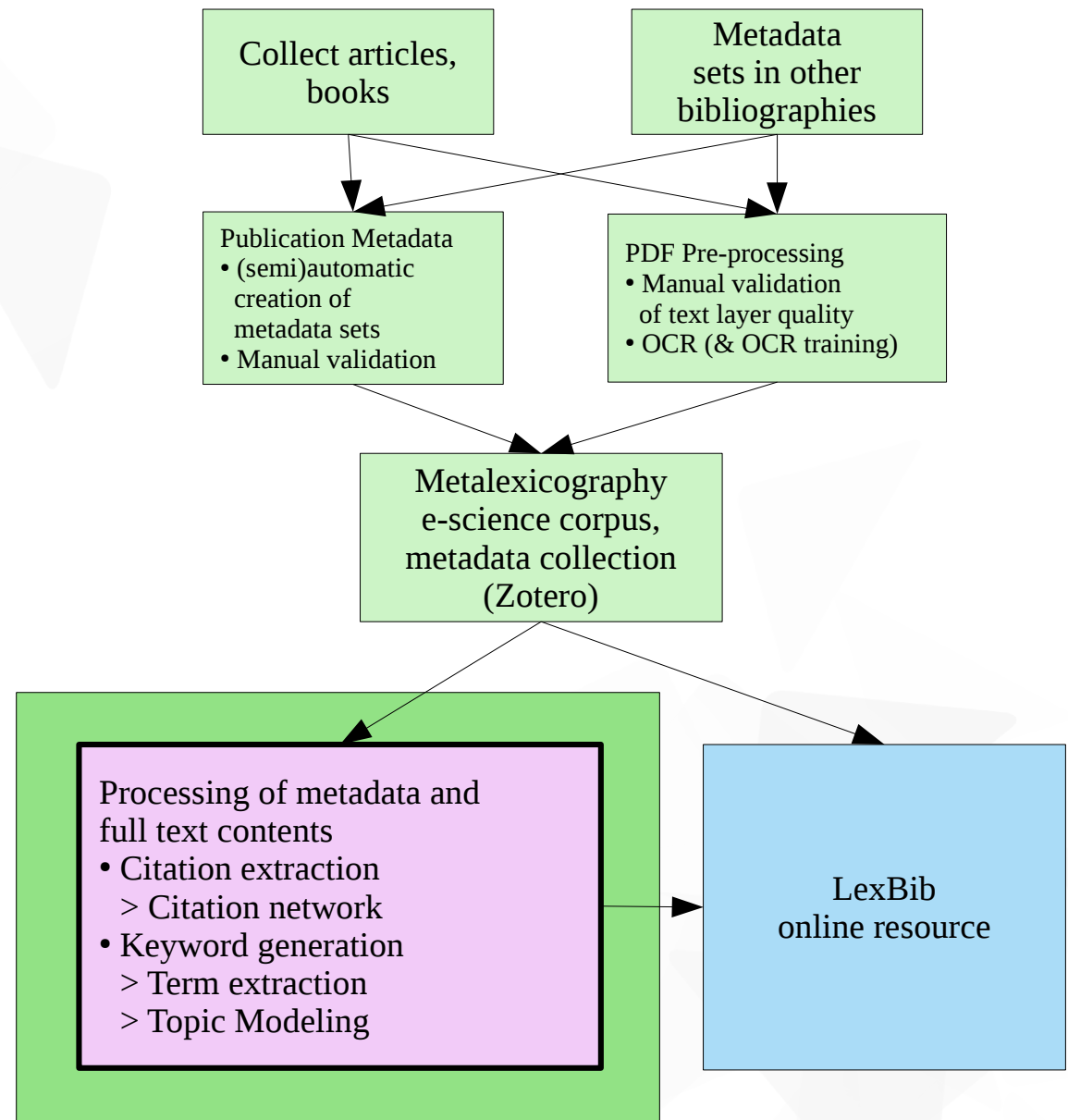
Green: e-science corpus: 2nd iteration

- ▼ Papers from other journals
- ▼ Papers from other conferences
 - ▼ Retrieval also by citation network created in 1st iteration
- ▼ Other languages: DE, ES, FR
 - ▼ Metadata, Full Texts
- ▼ Items from before 2000
 - ▼ PDF manipulation
 - ▼ OCR training
- ▼ Other item types
 - ▼ Books
 - ▼ Dissertations
 - ▼ Press articles
 - ▼ Blog posts, websites



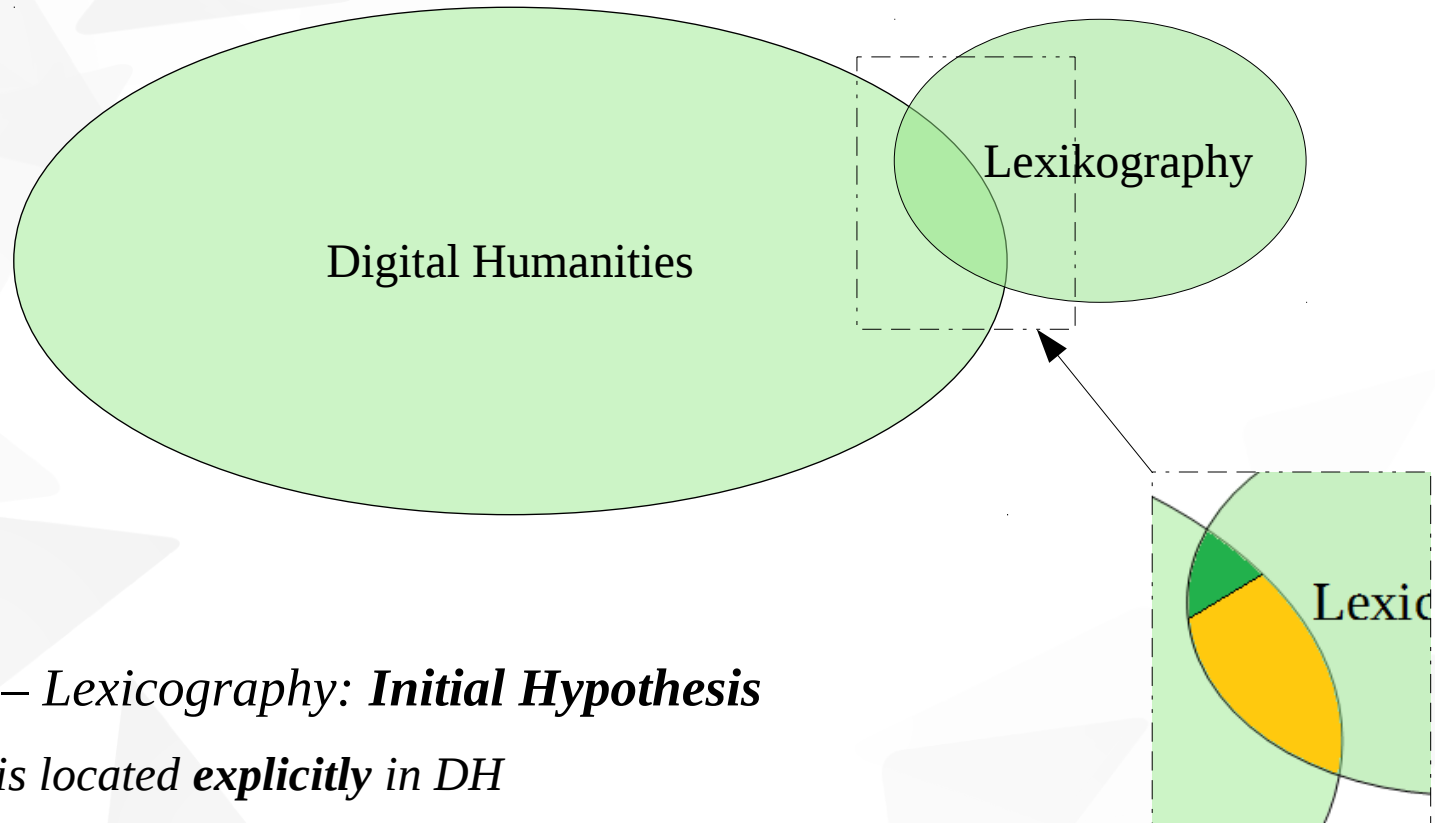
LexBib project: Module 2

- ▼ **Red:** Computational text analysis
- ▼ Citation extraction
 - ▼ PDF parsing: “GROBID”
 - ▼ Manual annotation of training material
 - ▼ Citation network
 - ▼ Item-Relations “citing“, “is_cited_by“
- ▼ Topic Modeling
 - ▼ LDA (“Mallet”)
- ▼ Term Extraction
 - ▼ “Trex“ (Schäfer et al.)



Preliminary Study: Lexicography and DH (DHd 2018)

Lindemann, Kliche & Kutzner 2018



- ▼ *Overlaps DH – Lexicography: **Initial Hypothesis***
- ▼ ***Only a part** is located **explicitly** in DH*
 - ▼ *...and publishes in a DH-journal or at a DH conference*
 - ▼ *...and describes own research as DH.*
- ▼ ***The bigger part** may be located **implicitly** in DH.*

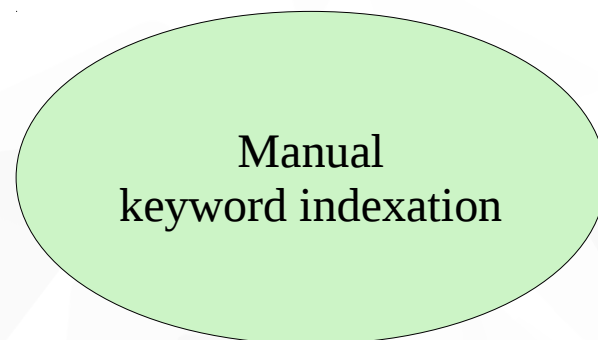
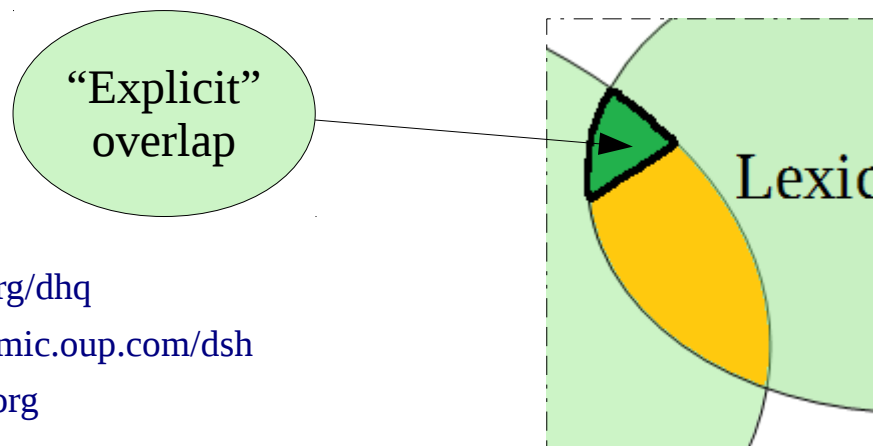
Manual search: Lexicography @ DH

Resources:

- Digital Humanities Quarterly: <http://www.digitalhumanities.org/dhq>
- Digital Scholarship in the Humanities (ex LLC): <https://academic.oup.com/dsh>
- TEI Journal of the Text Encoding Initiative: <http://jtei.revues.org>
- DHCommons: <http://dhcommons.org>
- DH Conferences: <http://digitalhumanities.org/dh-abstracts/search> | <https://dh2017.adho.org/program/abstracts/>

Keywords of 31 relevant contributions:

- Text Encoding / Markup / Encoding Formats / XML 31
- E-dictionaries / Visualization of lexical data 14
- Historical Lexicography 10
- Corpus Linguistics 9
- Linking Lexical Resources 4
- NLP-Lexicon 2
- Bilingual Dictionary Drafting 2
- Author dictionaries 2
- Dialectology 1



Lexikography, DH e-science corpora

▼ Definition of overlaps using DH methods

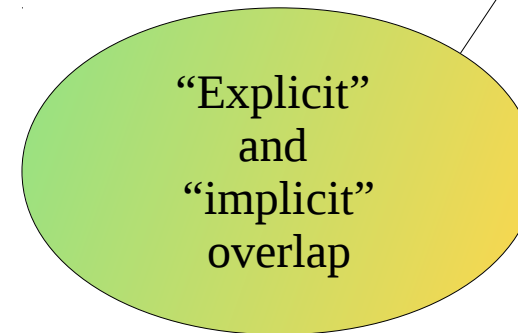
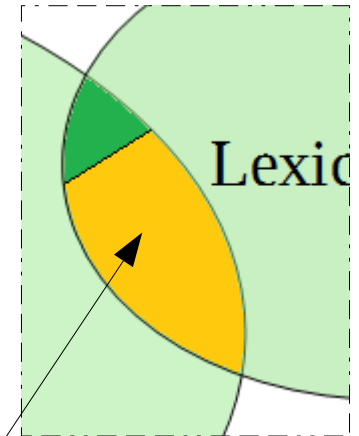
▼ two e-science corpora

▼ (1) DH corpus

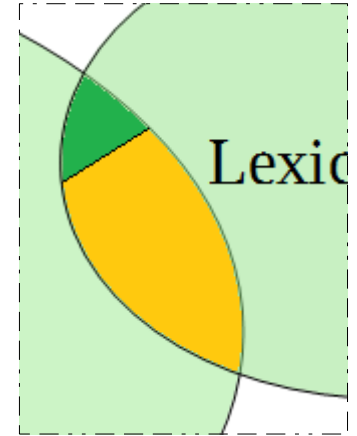
- ▼ DH journal articles
- ▼ DH conference papers
- ▼ DH handbook articles

▼ (2) Lexicography corpus

- ▼ Lexicography journal articles
- ▼ Lexicography conference papers
- ▼ Lexicography handbook articles



Lexicography corpus v1 (Spring 2018)



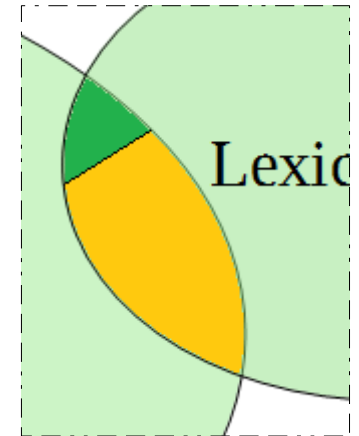
Lexicography: e-science corpus

▼ Proceedings of Euralex	▼ 782
▼ Lexikos Journal	▼ 376
▼ International Journal of Lexicography	▼ 282
▼ Journal of the Dictionary Society of North America (DSNA)	▼ 257
▼ Proceedings of eLex	▼ 202
▼ HSK 5/4: Dictionaries: An International Encyclopedia of Lexicography (Gouws et al., 2013)	▼ 110
▼ Routledge Handbook of Lexicography (Fuertes-Olivera, 2018)	▼ 47

DH corpus v1 (2018)

Digital Humanities e-science corpus

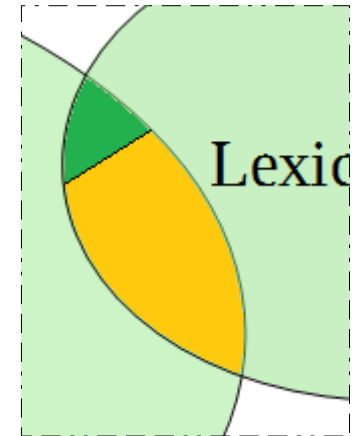
- ▼ Digital Scholarship in the Humanities (DSH/LLC) ▼ 886
- ▼ Digital Humanities Quarterly (DHQ) ▼ 284
- ▼ Digital Studies/Le champ numérique ▼ 152
- ▼ Journal of the Text Encoding Initiative ▼ 63
- ▼ Blackwell Companion to Digital Humanities (Siemens et al. 2004) ▼ 37



Lex_DH corpus v1 (2018)

Full text corpora

- ▼ Digital Humanities Corpus 1,380 Full texts
- ▼ Lexicography Corpus 1,919 Full texts (out of 2,057)
- ▼ Total 3,299 Full texts



Text Analytics: 3+ Methods

▼ Topic Modeling:

- ▼ LDA: Mallet
- ▼ Lemmatized Text
- ▼ Suppression of stop words

McCallum 2002

▼ Term Extraction:

- ▼ „TrEx“: Term extraction for CWB corpora
 - ▼ NN, NN-NN and NN-NN-NN patterns
 - ▼ Weirdness ratio: $\frac{\text{Relative Frequency in Lex_DH Corpus}}{\text{Relative Frequency in Ref.-Corpus}}$
 - ▼ Reference corpus: BNC

Schäfer et al. 2015

▼ Citation Analysis:

▼ GROBID:

- ▼ “GeneRation Of Bibliographic Data”
- ▼ Input: PDF full texts
- ▼ Output: TEI-XML

Romary & Lopez 2015

▼ Term Extraction:

- ▼ GROBID

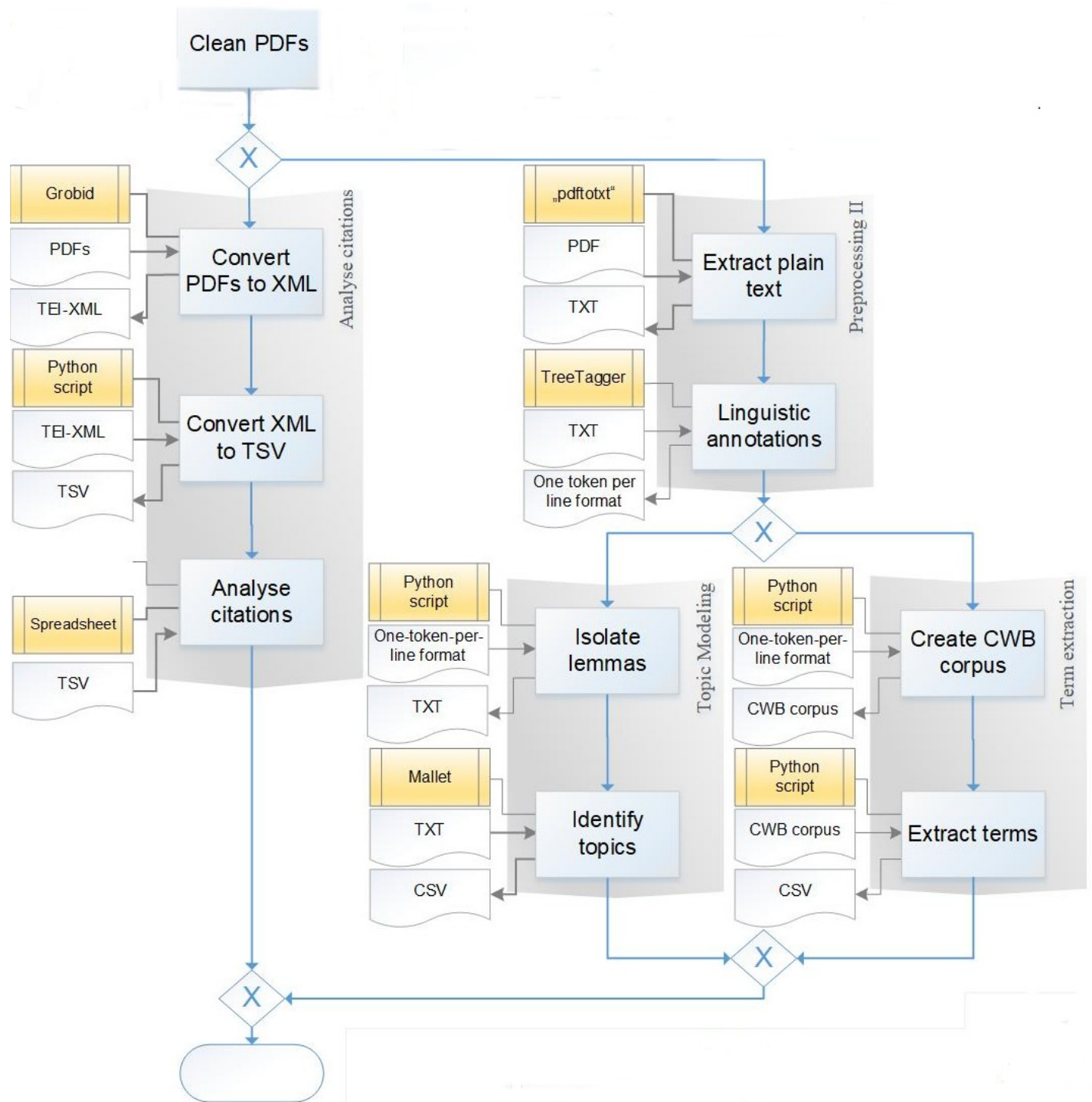
Lopez & Romary 2010

▼ Argumentative Zoning

▼ Citation Function

Teufel et al. 2002, 2006

Lex_DH e-science corpus text analytics workflow model



Citations

Normalization

Author = Last name of first author,
lowercase, [a-z]
Title = lowercase, [a-z]

Validation category 1

Triple **Author, Year, Title**
found in Lex_DH corpus as item

Validation category 2

Triple **Author, Year, Title**
found in Lex_DH corpus as item, and:

- $Lev(Author_i, Author_j) \leq 2$
- $\Delta(Year_i, Year_j) \leq 1$
- $Lev(Titel_i, Titel_j) \leq 8$

```
wiegand 2013 textualstructuresinprinteddictionariesan
wiegand 2013 textualstructuresinprinteddictionaries

wiegand 2013 macrostructuresinprinteddictionariesgouw
wiegand 2013 microstructuresinprinteddictionaries

kypridemou 2013 narrativesimilarityascommonsummarypaperp
kypridemou 2014 narrativesimilarityascommonsummary

mccarty 2005 treeturfcentrearhipelagoorwildacremetap
mccarty 2006 treeturfcentrearhipelagoorwildacre

siepmann 2006 collocationcolligationandencodingdiction
siepmann 2005 collocationcolligationandencodingdictionaries

orli 2006 theoxforddzscomprehensiveenglishslovenia
orli 2006 theoxforddzscomprehensiveenglishslovenian

bergenholtz 2003 userorientedunderstandingofdescriptivepr
bergenholtz 2003 userorientedunderstandingofdescriptive

fuertesolivera 2016 acambrianexplosioninlexicographysomerefl
fuertesolivera 2016 acambrianexplosioninlexicography
```

Result counts:

val_cat 1: 760
val_cat 2: 1639

Lex_DH Citation Network

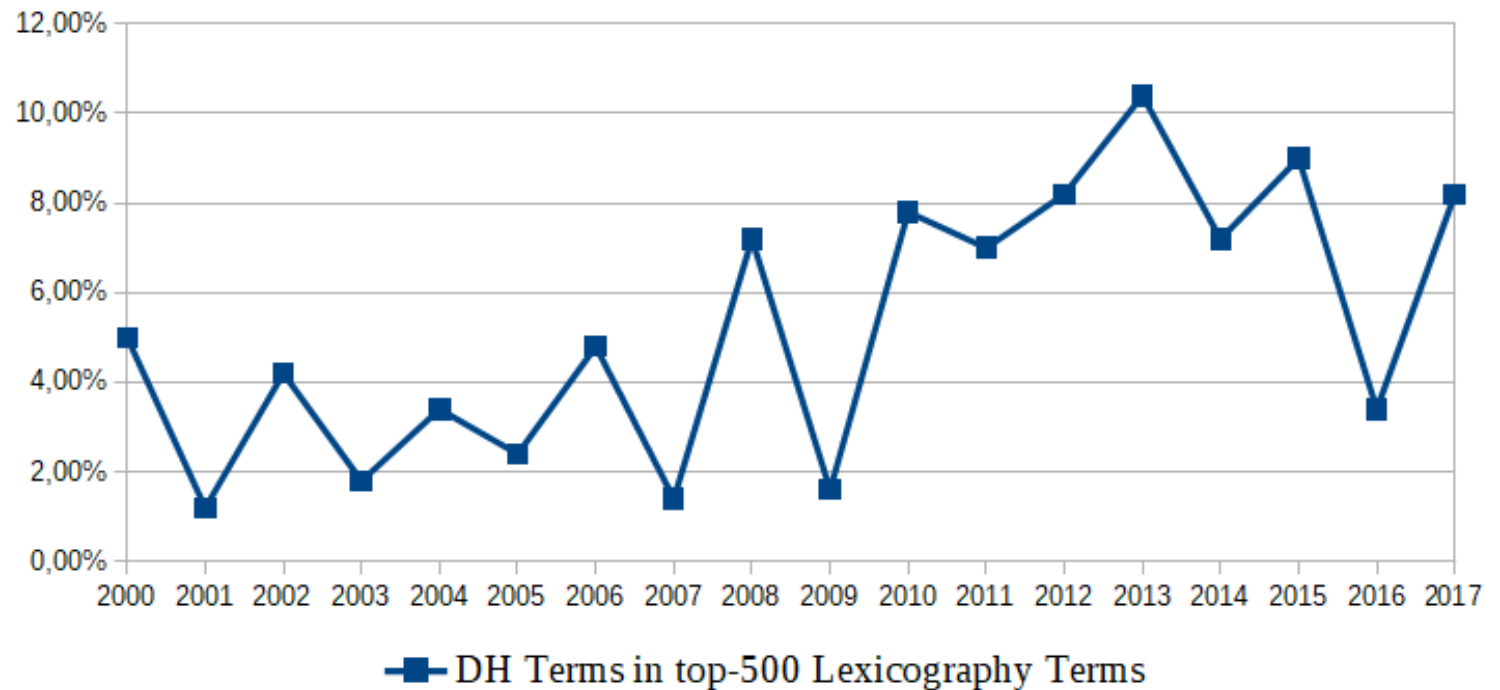
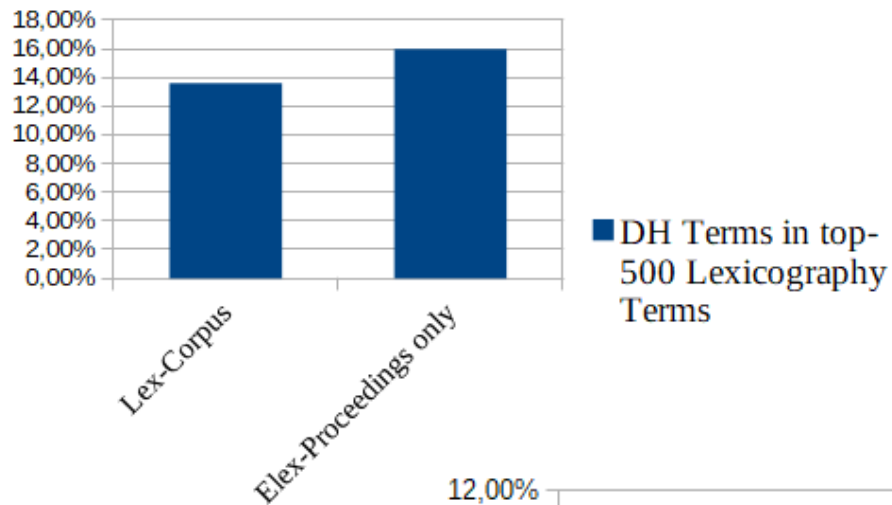
1,108	citing articles	i.e. articles in Lex_DH corpus that cite one or more articles in Lex_DH corpus
	709 Lex 399 DH	
2,432	citation pairs	i.e. citations in Lex_DH corpus that point to another Lex_DH article
	1,674 Lex 758 DH	
1,057	cited articles	i.e. articles in Lex-DH corpus that are cited in one or more Lex_DH articles
	715 Lex 342 DH	

742	Pairs	DH → DH
16	Pairs	DH → Lex
1,659	Pairs	Lex → Lex
15	Pairs	Lex → DH

Results: Term Extraction Lex_DH

DHTerms sortHits	DHTerms sortWeird	LexTerms sortHits	LexTerms sortWeird	Top LexTerms found in DH
markup	website	lexicography	dictionary article	website
internet	pdf	lemma	access structure	lemmatization
website	xml	internet	dictionary user	wordnet
authorship attribution	stemma	dictionary use	lemma sign	reference corpus
metadata	text mining	macrostructure	text reception	corpus query
digitization	authorship attribution	polysemy	multiword	search engine
pdf	blog	multiword	website	internet
modeling	text classification	dictionary article	dictionary consultation	web site
text analysis	cyberinfrastructure	multi-word	word formation	web page
xml	search engine	access structure	lexicography	text box
stemma	feature selection	source language	corpus evidence	print version
blog	url	text production	text production	subcorpus
stylometry	classification accuracy	website	lemmatization	frequency list
cyberinfrastructure	web page	word formation	dictionary research	crowdsourcing
url	web site	dictionary user	function theory	web interface
dh	php	lemmatisation	article stretch	corpus research
search engine	crowdsourcing	lemma sign	word sketch	language documentation
avatar	open-source	lemmatization	wordnet	word alignment
curation	base text	text reception	reference corpus	hyperlink
part-of-speech	internet	wordnet	translation equivalent	wikipedia
intertextuality	text categorization	translation equivalent	dictionary information	search interface
text mining	test text	pdf	pdf	blog
word frequency	text reuse	definition	lemma list	source word
open-source	book history	target user	dictionary making	text genre
php	metadata	dictionary consultation	definition	word sense disambiguation

Results: Term Extraction Lex_DH



Results: Topic Modeling Lex_DH



david.lindemann@uni-hildesheim.de
fritz.kliche@uni-hildesheim.de
heid@uni-hildesheim.de

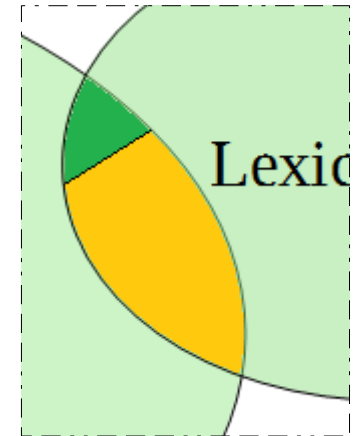
#	DHLEX									
25	100	0	digital library collection research information datum http project access resource material www metadata web record provide document object user archive							
27	100	0	digital medium text reading work literary literature technology press form space ai							
46	100	100	0	digital library collection research information	scholarly sci					http history
15	99	100	0	digital medium text reading work literary litera	tion editor					
49	99	100	0	digital medium text reading work literary litera	er charac					
8	95	100	0	digital humanity humanities computing resear	er analysis color page panel colour shape fig letter show document					
23	95	99	1	text edition manuscript textual electronic digit	akespeares count passage poet test rhyme work block figure length					
16	90	99	1	text word author feature analysis authorship r	del discourse time code summary avatar world video adventure interactive					
41	89	99	1	text word author feature analysis authorship r	object argument pattern phrase structure type framenet role class information lexicon					
40	87	95	5	image visual film visualization figure sound m	am group develop development member education technology online time provide make participant					
13	84	95	5	shakespeare word poem line play text poetry	on markup type initiative http structure issue ontology attribute journal datum					
44	84	95	5	shakespeare word poem line play text poetry	le information citation study relationship datum twitter represent communication edge show					
4	81	90	10	game narrative event story player character a	ness variation medieval stemma tree scribe figure deaf version reading distance					
24	64	90	10	game narrative event story player character a	er analysis color page panel colour shape fig letter show document					
10	60	40	card table number page figure total show high compare study sample result journal list average percentage type distribution author score							
31	56	44	corpus word text language frequency english linguistics datum study linguistic analysis speech list discourse include speak research table genre million							
5	54	46	card case make give point present approach information study work form type part fact analysis question find provide result structure							
32	48	52	database system tool text software datum project file computer user application design work interface create code figure process develop xml							
28	44	56	card word language corpus method result system text proceeding sentence computational evaluation approach error automatic algorithm number extract tag datum							
36	44	56	legal law copyright business term court political property company acronym register public act trade government state land market financial fair							
3	33	67	woman slang gender female male man term sex fem	64	36	dialect map datum linguistic				
42	33	67	card book work history page author write year publis	60	40	card table number page figure				
20	31	69	chinese arabic irish english character china canadia	56	44	corpus word text language f				
39	27	73	word form rule morphological noun lexicon language	54	46	card case make give point p				
7	26	74	word make people geen thing time find good sense	48	52	database system tool text sc				
17	24	76	english johnson london language john johnsons ed	44	56	card word language corpus				
34	24	76	estonian adv language hindi arabic indian proper in	44	56	legal law copyright business				
43	23	77	proverb japanese boyer french page royal expressio	33	67	woman slang gender female				
18	20	80	spanish card del diccionario los una madrid por cata	33	67	card book work history page				
37	20	80	gabonese gabon french stein part polis des work ma	31	69	chinese arabic irish english				
38	18	82	des french les card dan dictionnaire une qui langue							
6	15	85	italian della dizionario del anglicism lingua italiana ci							
12	14	86	translation equivalent language bilingual translate e							
14	12	88	dutch van language dictionary slovene frisian datab							
21	11	89	der german die und card das von des den mit ein zur							
48	10	90	user search online card access electronic information http figure www web internet datum resource result link tool content query provide							
26	5	95	idiom polish czech expression phraseological idiomatic usage syn maiv idioms variant slownik mienh light unit variation maaih praha figurative nyei							
2	4	96	semantic lexical sense card meaning word relation definition language concept wordnet synonym metaphor category lexicon semantics conceptual thesaurus cognitive structure							
30	4	96	english oed quotation dictionary word citation oxford editor murray evidence source historical middle etymology supplement entry language press med american							
45	3	97	article card item structure dictionary text access wiegand give address type cross-reference outer microstructure form lemma reference follow partial relation							
1	2	98	swedish language danish dictionary lemma word norwegian project romanian list form card editor version croatian lexicon publish compound nordic german							
29	2	98	card dictionary learner word student study english language test vocabulary information learners participant							
0	1	99	language card dictionary african south shona africa speaker english ndebele community lexicography cultura	0	100	term terminology concept domain termino				
47	1	99	dictionary user card text datum lexicography information lexicographic function situation tarp bergenholtz the	0	100	dictionary word entry card english definiti				ific artic
9	0	100	die van word wat nie vir afrikaans woordeboek meet hierdie dit woordeboeke card dat heat aan kan gebruik	0	100	card dictionary lexicography language in				nce sci
11	0	100	term terminology concept domain terminological field card information language specialized definition knowle	0	100	dictionary word sotho class noun stem le				e electr
19	0	100	dictionary word entry card english definition include meaning language sense headword information label list	0	100	collocation card corpus collocate lexical				
22	0	100	card dictionary lexicography language international lexicographic linguistics corpus research dictionaries edit							
33	0	100	dictionary word sotho class noun stem lemma form northern african schryver language zulu user prinsloo ver							
35	0	100	collocation card corpus collocate lexical verb word noun ofthe pattern combination collocational expression sketch adjective unit phrase meaning base semantic							

Results: Topic Modeling



Preliminary Study Lex/DH: Results

- ▼ Contribution to mutual consideration of research communities
 - ▼ Citation overlap (“explicit overlap“) smaller than “implicitly“ overlapping Topics
- ▼ Overlapping Topics
 - ▼ Common Methods
 - ▼ Examples: Corpus Linguistics; Language Documentation
 - ▼ Common “third“ research interests
 - ▼ Examples: Gender; juridical issues (copyright etc.)
- ▼ Outlook
 - ▼ Close reading in overlapping / disjunctive sets
 - ▼ If terms / topics overlap: Which texts are we talking about?
 - ▼ If citations overlap, but extracted terms / topics do not: Are we talking about the same, using different words?



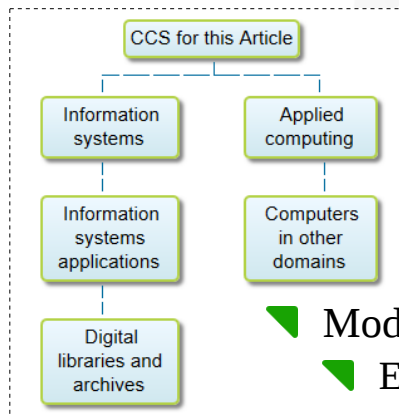
Term extraction tests

- ▼ Example: Top 20 terms, extracted from a single full text (here, a 2014 Euralex keynote speech)
 - ▼ Reference corpus BNC
 - ▼ Reference corpus LexBib
- ▼ Mapping to existing keyword lists
 - ▼ Obelex-meta
 - ▼ as silver standard for evaluation
 - ▼ as first basic grid for keyword ontology
 - ▼ [etc.]

<i>Top 20 Terms (Ref. BNC)</i>	<i>Top 20 Terms (Ref. LexBib)</i>
text reception	information-on-demand
text production	on-demand
dictionary function	data repository
multiword	user friendliness
word formation	user orientation
production dictionary	text reception
user orientation	production dictionary
information-on-demand	text production
data repository	dictionary function
dictionary entry	repository
user friendliness	valency
internet	guidance
markup	orientation
on-demand	concord
valency	word formation
concord	scenario
dictionary	markup
corpus	production
collocation	classification
language processing	advance

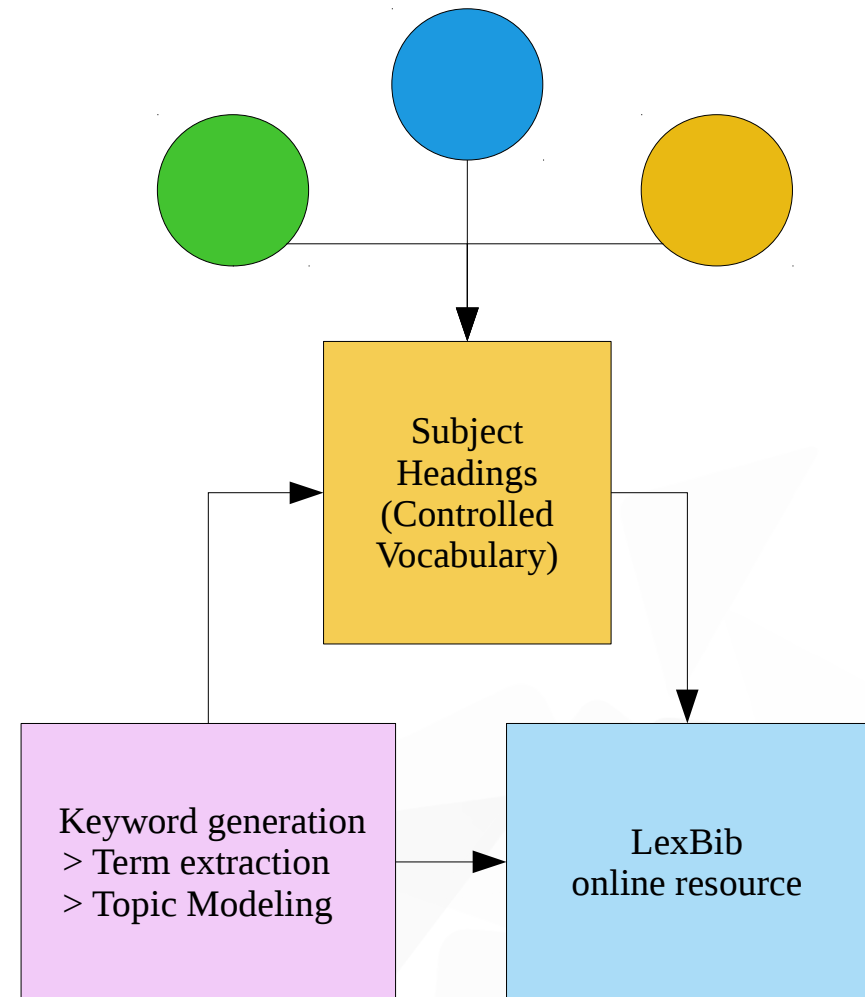
Additional Module: Keyword / Subject heading ontology

- Development and application of an ontological controlled vocabulary for subject headings
- Sources:
 - e-science corpus: extracted terms
 - e-science corpus: author keywords
 - Obelex-meta keywords
 - Wiegand's concept taxonomies
 - Library of Congress Subject Headings
 - GND (Gemeinsame Normdatei)



Model: ACM CCS

Example Bamman & Crane 2008 @ ACM

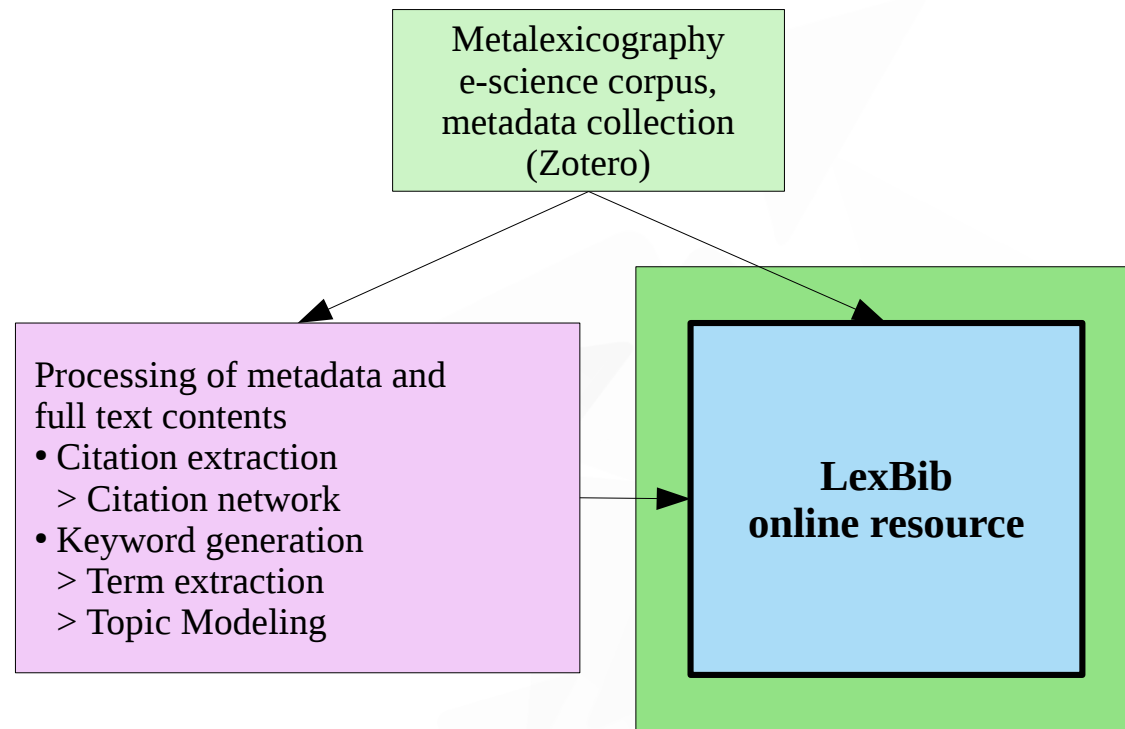


LexBib project: Module 3

▼ Blue: Creation of online bibliography

▼ Features:

- ▼ Enhanced search functions
- ▼ Links to full text download pages
- ▼ Persons / events: Profile pages
- ▼ API for metadata, BibTeX, RDF
- ▼ is-container-of / is-contained-in
- ▼ Browsible citation network
- ▼ is-review-of / is-reviewed-by
- ▼ Keywords: Results of red module
- ▼ Browsible keyword ontology



LexBib project: Partners

Green

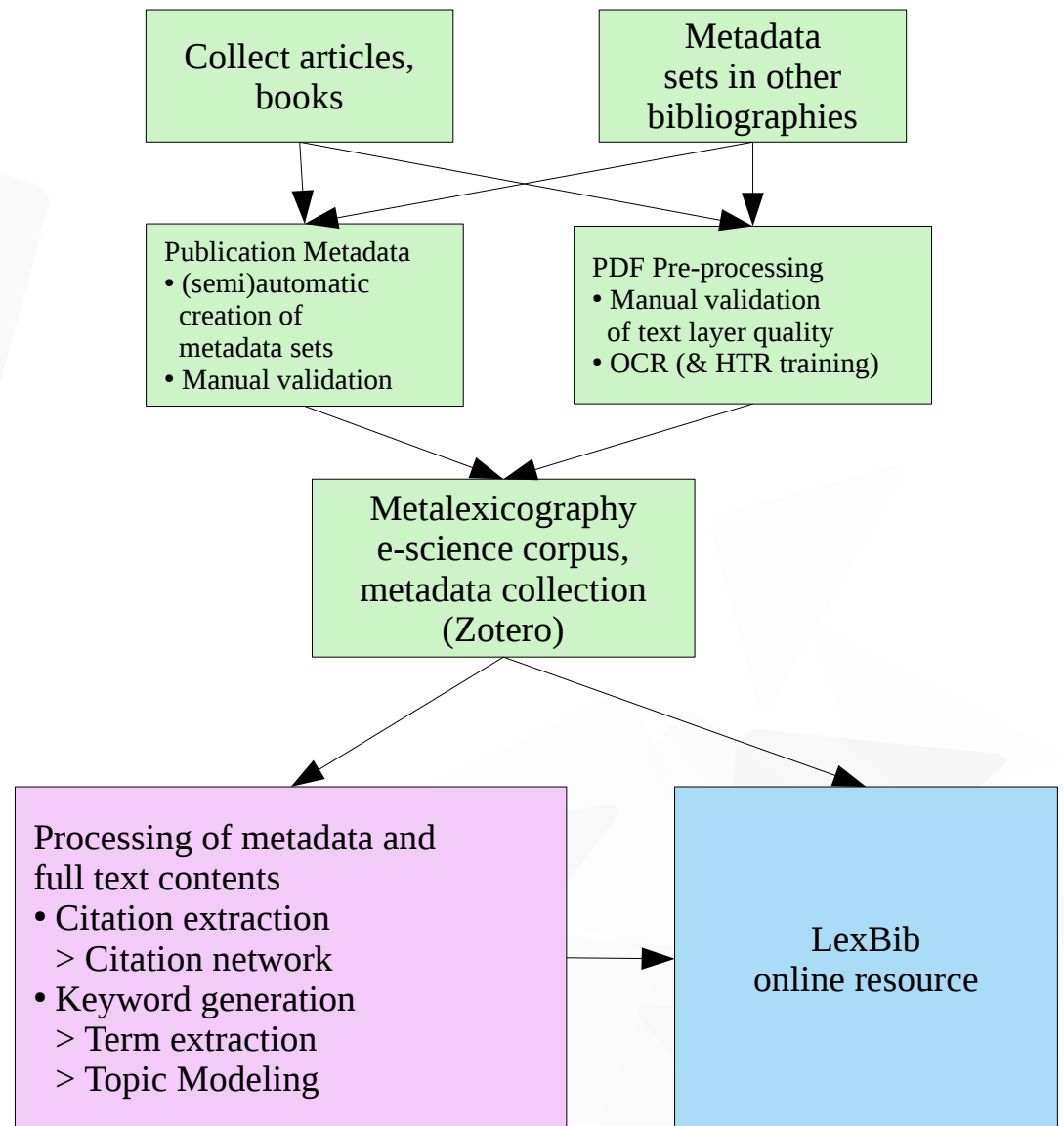
- ▼ Elexis
- ▼ Euralex
- ▼ IDS
- ▼ GLex UDC
- ▼ EMLex
- ▼ de Gruyter

Red

- ▼ Romary, Khemakhem (Inria, CMB)

Blue

- ▼ Clarin-D (IDS)
- ▼ Euralex
- ▼ Elexis



LexBib project: Long-term perspective

▼ Complete collections (before 2020)

- ▼ various languages: EN, DE, ES, FR
 - ▼ publication metadata
 - ▼ full texts

▼ LexBib Online Resource

- ▼ Bibliography items with unique IDs
 - ▼ publications
 - ▼ persons
 - ▼ events, places
- ▼ Keyword indexation
 - ▼ automatically extracted
 - ▼ mapped to controlled vocabulary

▼ After 2020

- ▼ **Blue:** Long-term hosting options
 - ▼ CLARIN server @ IDS
 - ▼ ELEXIS

▼ **Green:** Curation team and/or community?

- ▼ Author pages
 - ▼ editing of own page?
- ▼ Keyword indexation of new items
 - ▼ editing of own items (as in ACM)?

▼ **Red:** Tools as services

- ▼ Application to e-science corpora of other domains

david.lindemann@uni-hildesheim.de
fritz.kliche@uni-hildesheim.de
heid@uni-hildesheim.de

<http://www.zotero.org/groups/lexbib/items>



David Lindemann
Fritz Kliche
Ulrich Heid

Thank you for your attention

See detailed bibliography in the proceedings