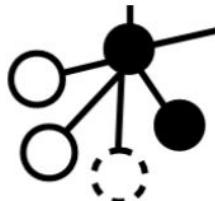


LOC-DB: A Linked Open Citation Database provided by Libraries. Motivation and Challenges.

*EXCITE Workshop 2017
March 30/31, 2017
GESIS, Cologne*



Tiessen, Jan (2007): Die Resultate im Blick?
ner/Döhler, Marian (Hrsg.): Agencies in W
Tondorf, Karin/Bahnmüller, Reinhard/Klages,
instrument. Anwendungspraxis, Probleme
sigma.

Touraine, Alain (1984): Le retour de l'acteur: e
Treiber, Hubert (1984): Warum man nicht die
Mikroskop den ganzen Elefanten zu sehen.



Overview

1. Motivation

Presenter: Kai Eckert, Stuttgart Media University

2. Infrastructure

Presenter: Anne Lauscher, Stuttgart Media University

3. Reference Extraction

Presenter: Akansha Bhardwaj, DFKI Kaiserslautern

Motivation

Scientists:

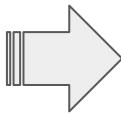
- ... search for **content**, not for books.
- ... need to check **citations**.
- ... need to know **who cites** the content.
- ... as a special case: need to know who cites **their** content.



Source: https://commons.wikimedia.org/wiki/File:Mad_scientist_transparent_background.svg

Libraries: What is in the catalog?

Basically: Things you can put on a shelf.



But there are databases!



Many commercial, incomplete databases.

Many smaller, focused databases.

Many library services, like the [OnLine Contents Fachausschnitte](#) (fka. OLC-SSG).

All are **incomplete** (most do not even try to be complete)

Most are **not freely available as Open Data**.

Most do **not contain citation references**.

What if...

... libraries would just do it?

- **Catalog** journals and collections (proceedings),
- with **all articles / chapters** as separate resources,
- with a structured form of all **citation references** in all articles,
- ideally with **links** to the cited resources.

Linked Open Citation Database

Goal: Answering the question:

What is needed (persons, resources, money)
for libraries to actually “do it”?

(i.e., not the creation of a complete and free citation database, at least for now...)

Method: Prototypical creation of a complete (in the sense of a subset) and free citation database with a focus on cataloging efficiency.

Consortium:



Deutsches
Forschungszentrum
für Künstliche
Intelligenz GmbH



HOCHSCHULE DER MEDIEN



Leibniz-Informationszentrum
Wirtschaft
Leibniz Information Centre
for Economics

Funded by:

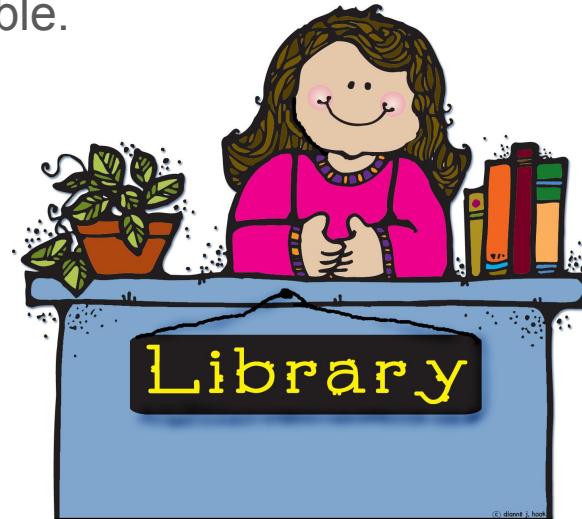


Focus on efficiency

- Reuse existing data (e.g., from publishers or other projects).
- Use OCR.
- Semi-automatic data extraction and linking.
- Streamline and automate the process wherever possible.
- Distributed database and cataloging process.

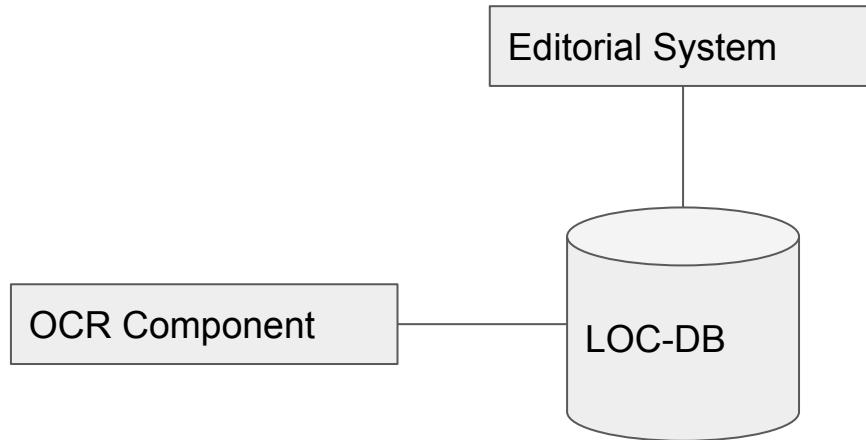
Desired answer: If X libraries use Y persons to do the LOC-DB cataloging, we manage to get Z percent of the content.

Hopefully with low X and Y and high Z ;-)



Infrastructure

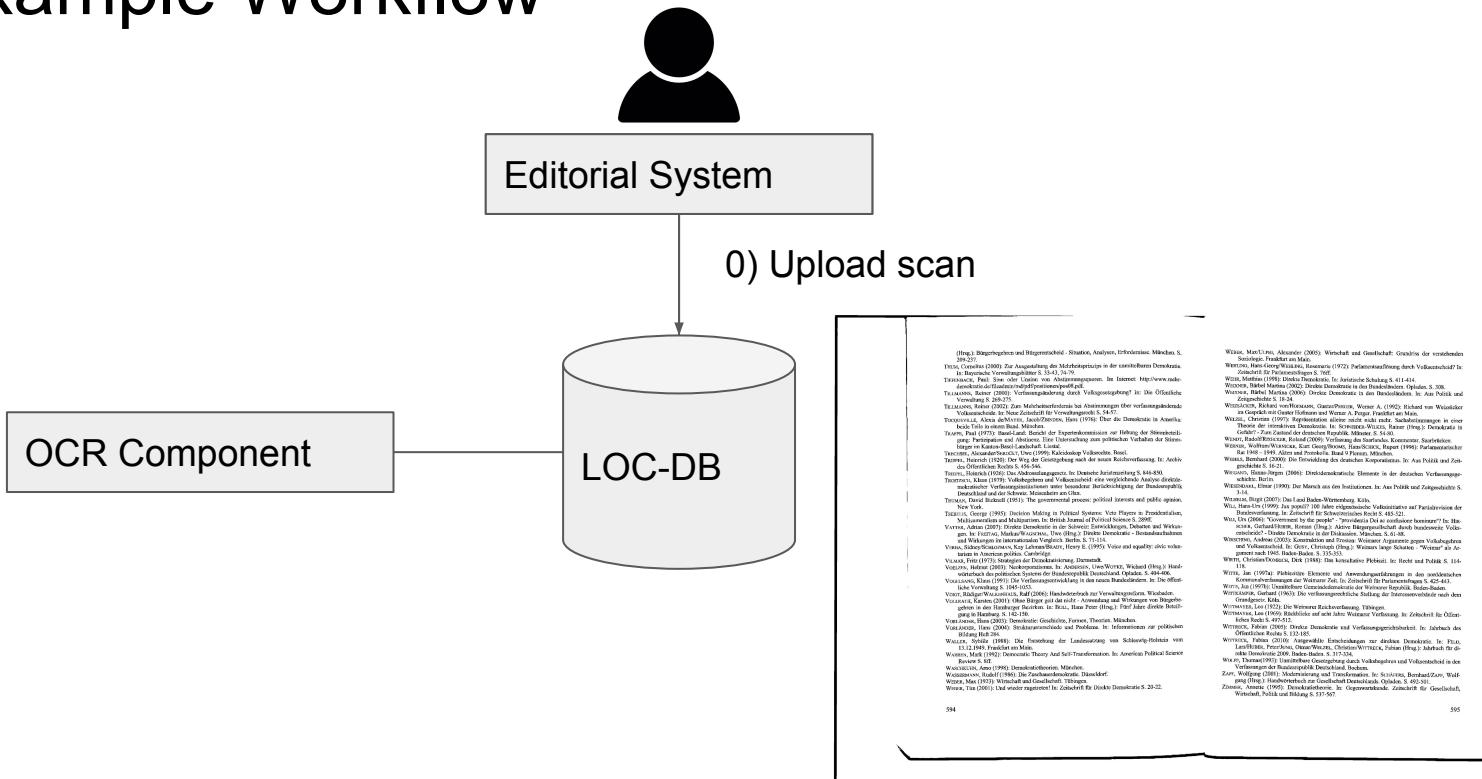
Overview



Main components:

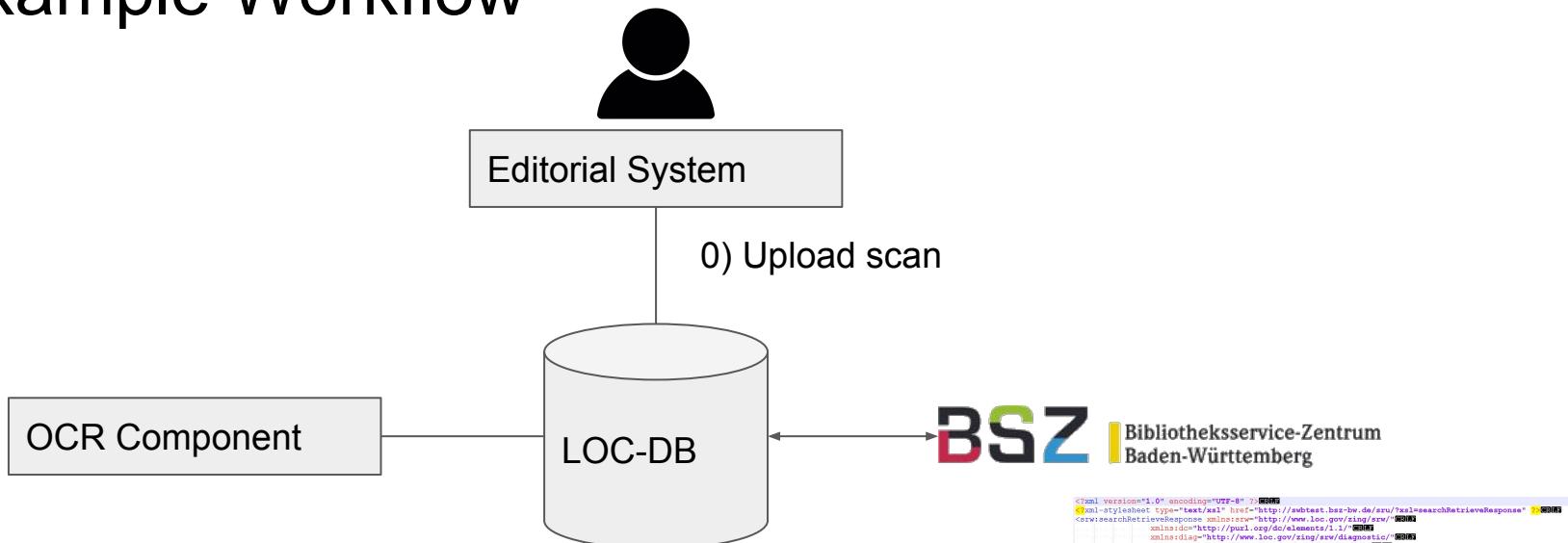
- Editorial System (GUI)
- Central Component (LOC-DB)
- OCR Component

Example Workflow



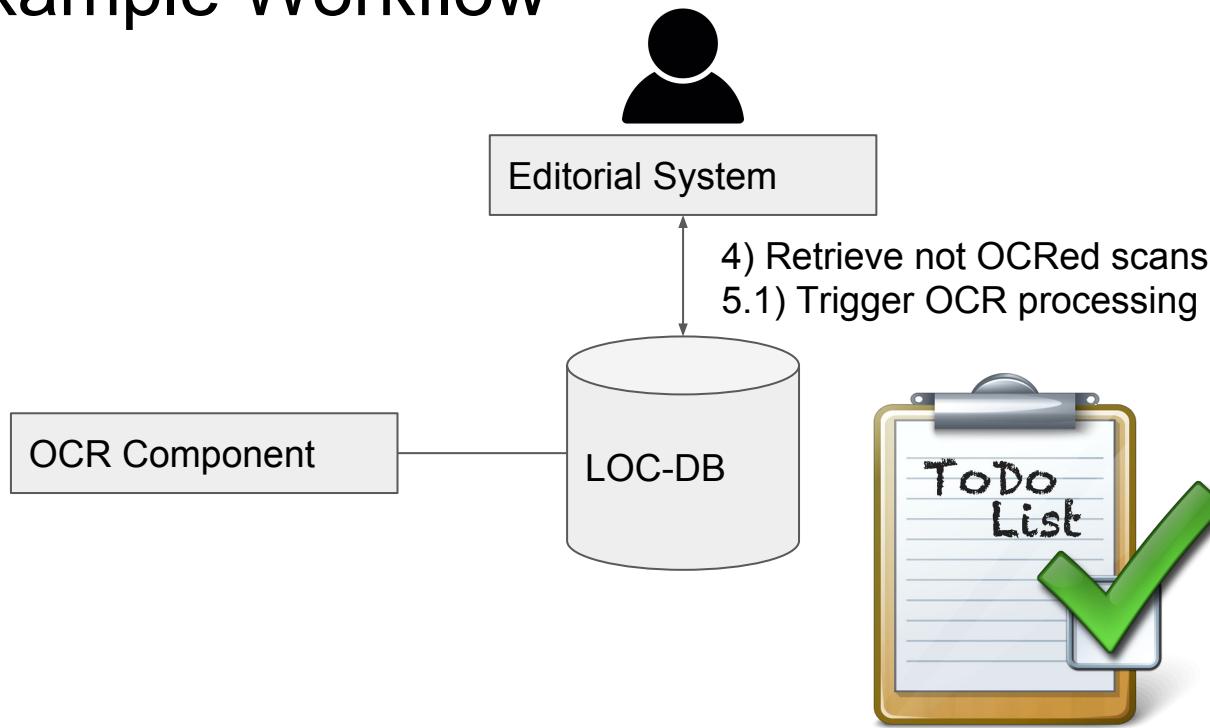
Sources of the images used: [1]

Example Workflow

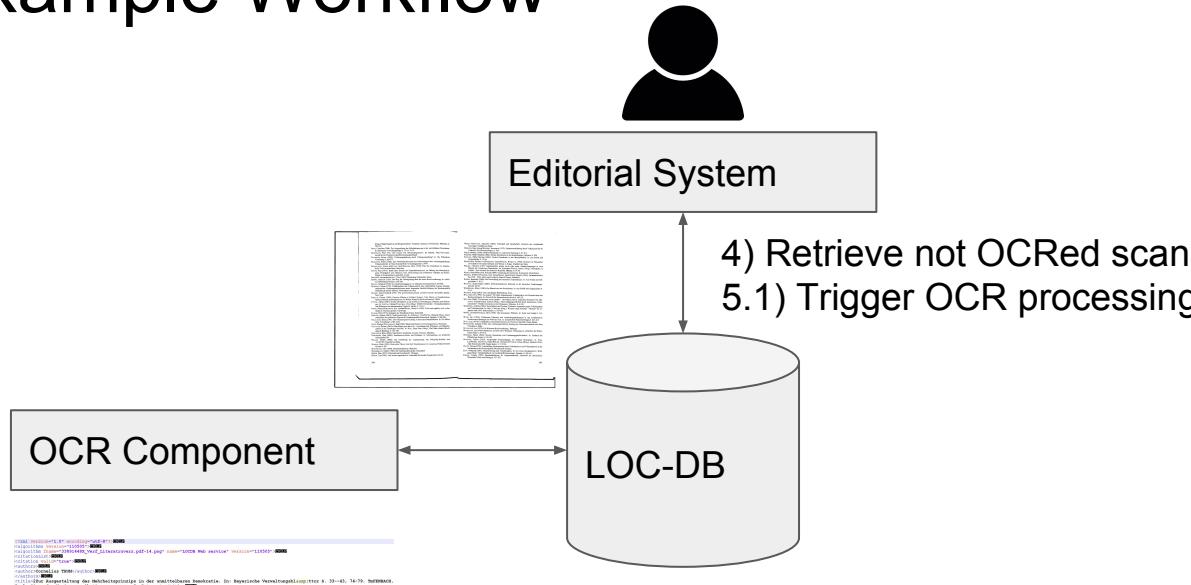


- 1) Save file
 - 2) Get related data from SWB-BSZ
 - 3) Create entry in the database

Example Workflow



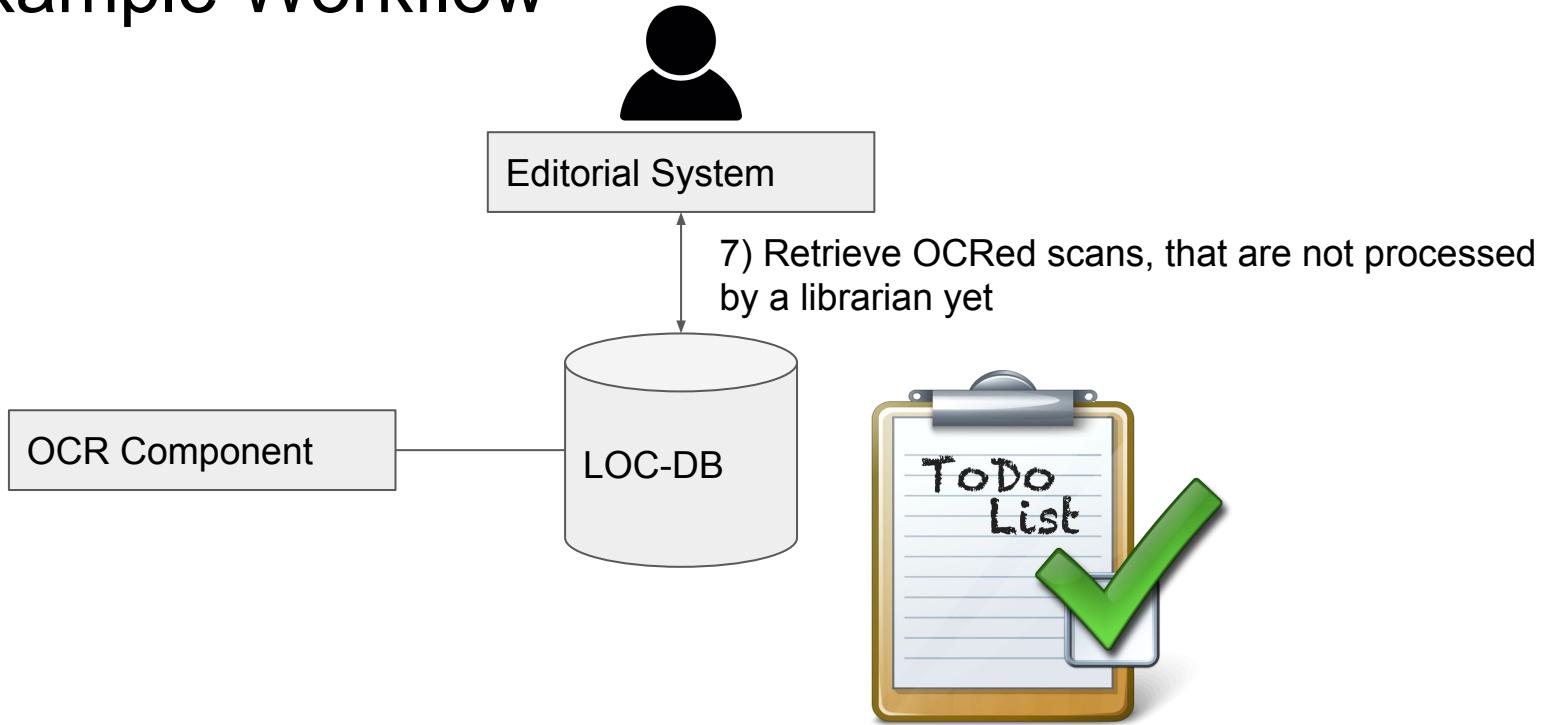
Example Workflow



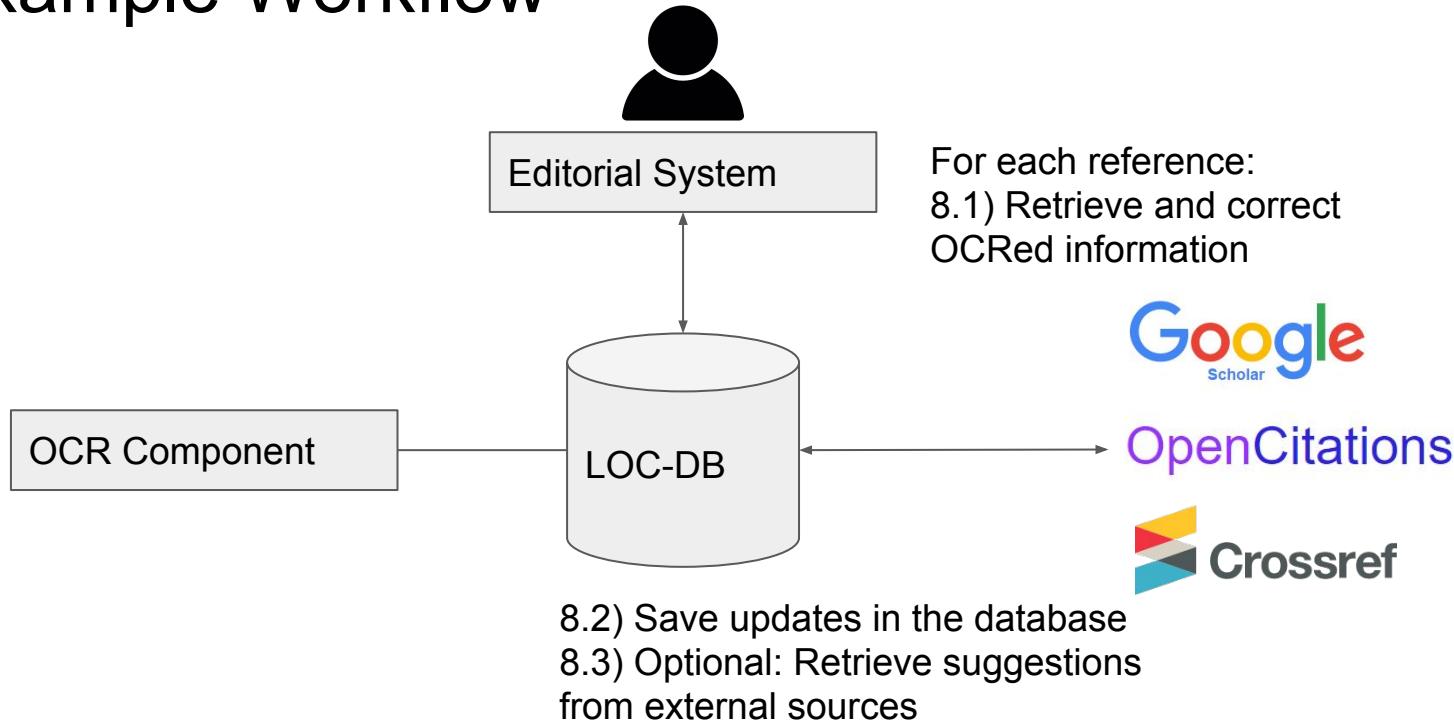
5.2) Trigger OCR processing

6) Retrieve OCRed data (e.g. coordinates of the references on the scanned page) and save it

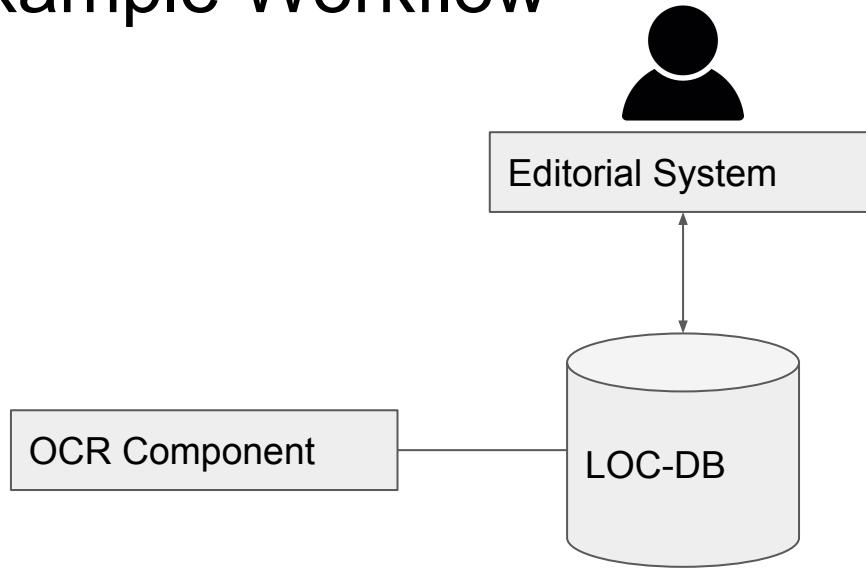
Example Workflow



Example Workflow



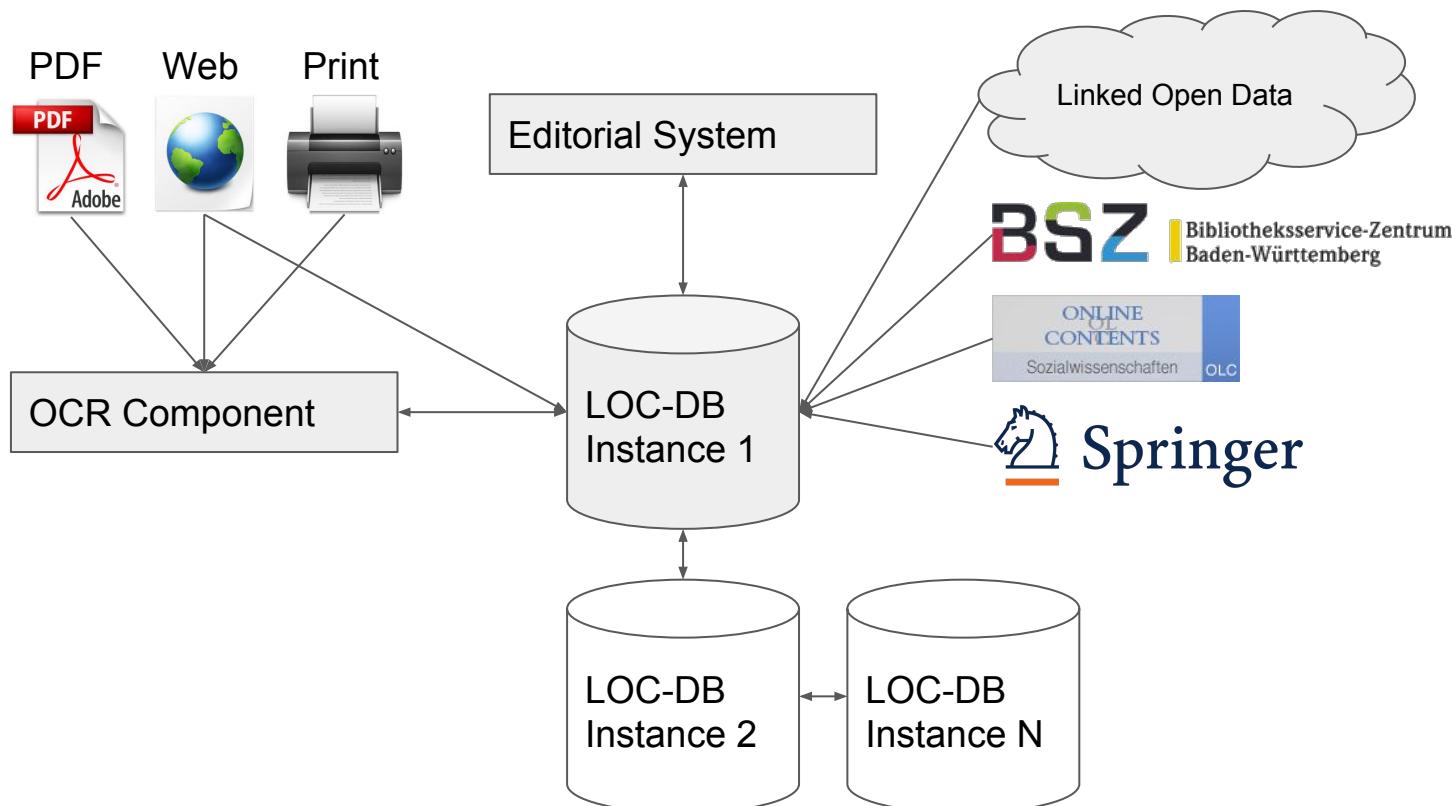
Example Workflow



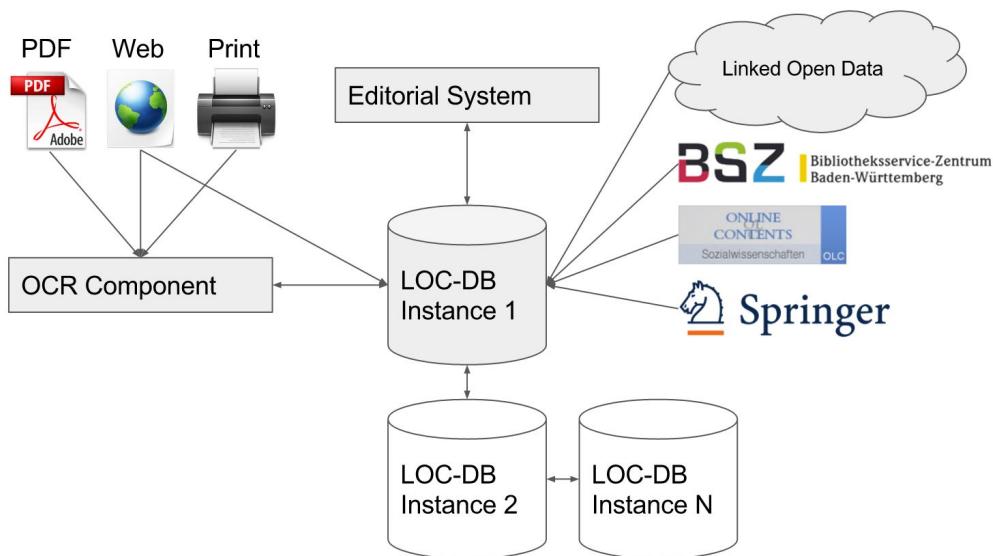
Result:

A linked open citation database extracted from reference lists
created by librarians

A bit more than just three components..



Technologies

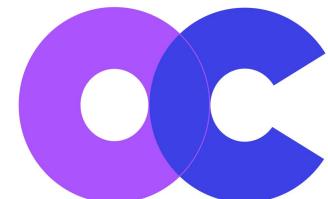


Main components:

- Editorial System (GUI)
 - [Angular.io](#)
 - [TypeScript](#)
 - [Bootstrap](#)
- Central Component (LOC-DB)
 - [Swagger](#)
 - [Node.js](#)
 - [MongoDB](#)
- OCR Component

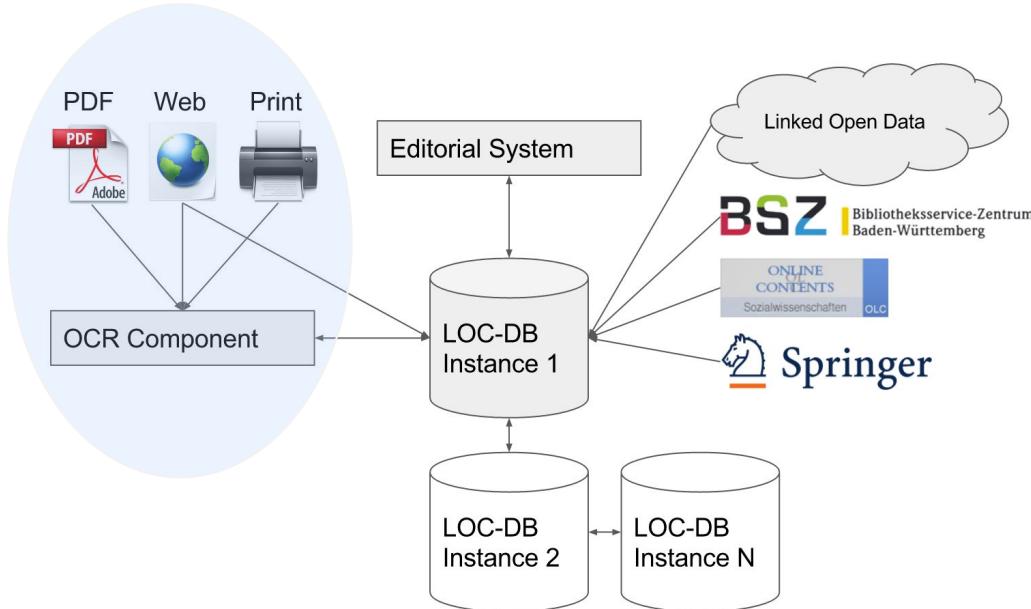
Data Model

- Inspired by the [OpenCitations](#) Data Model
- Extensions to support the management of scans and detected references in different status
- Example: BibliographicEntry (“a single reference”)
 - Identifier of the corresponding scan
 - Status
 - OCRed information, e.g. coordinates of the reference on the scan, title etc.



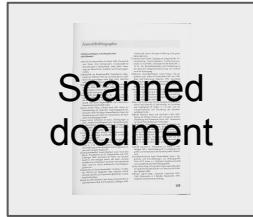
Reference Extraction

Overview



Main components:

- Editorial System (GUI)
- Central Component (LOC-DB)
- **OCR Component**



Scanned
document



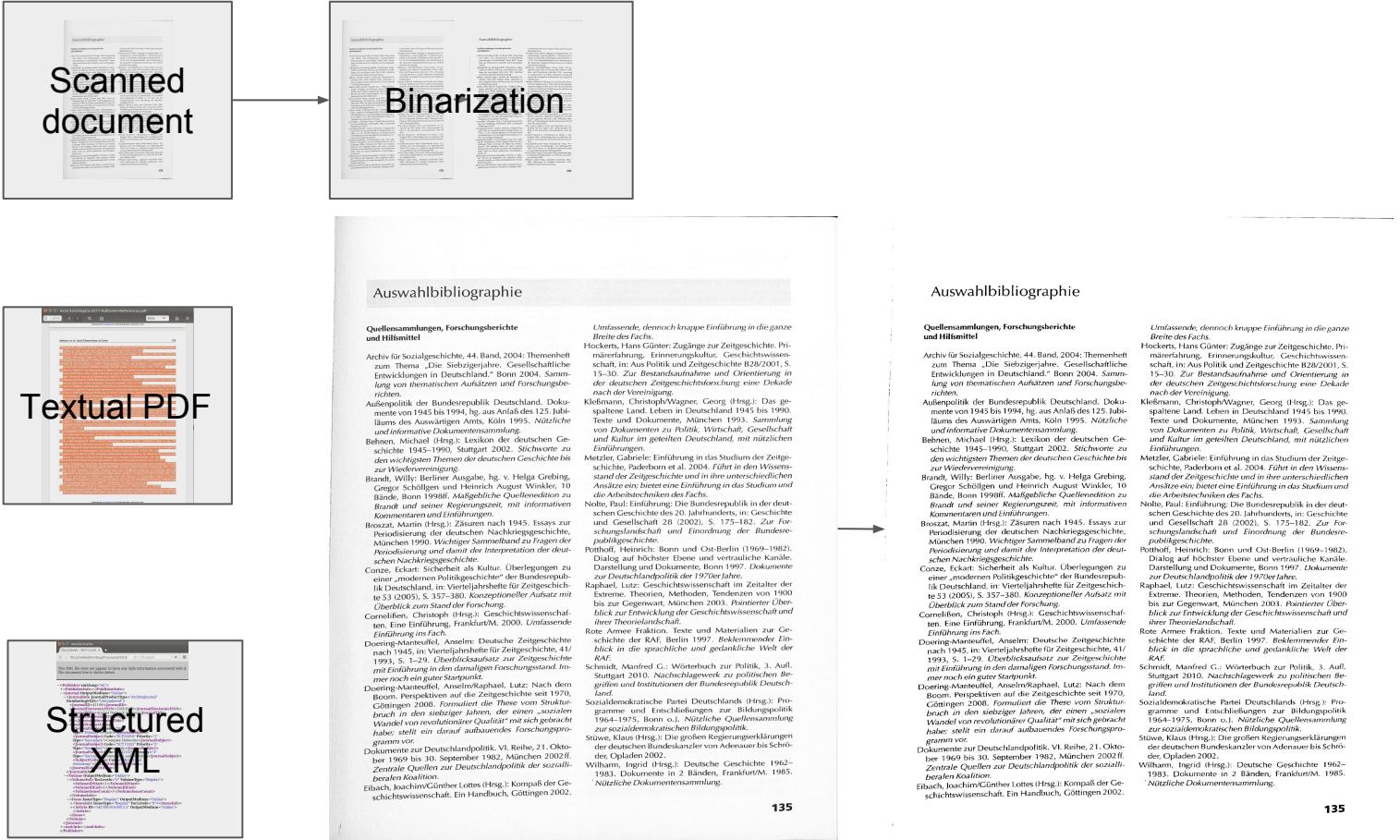
Textual PDF

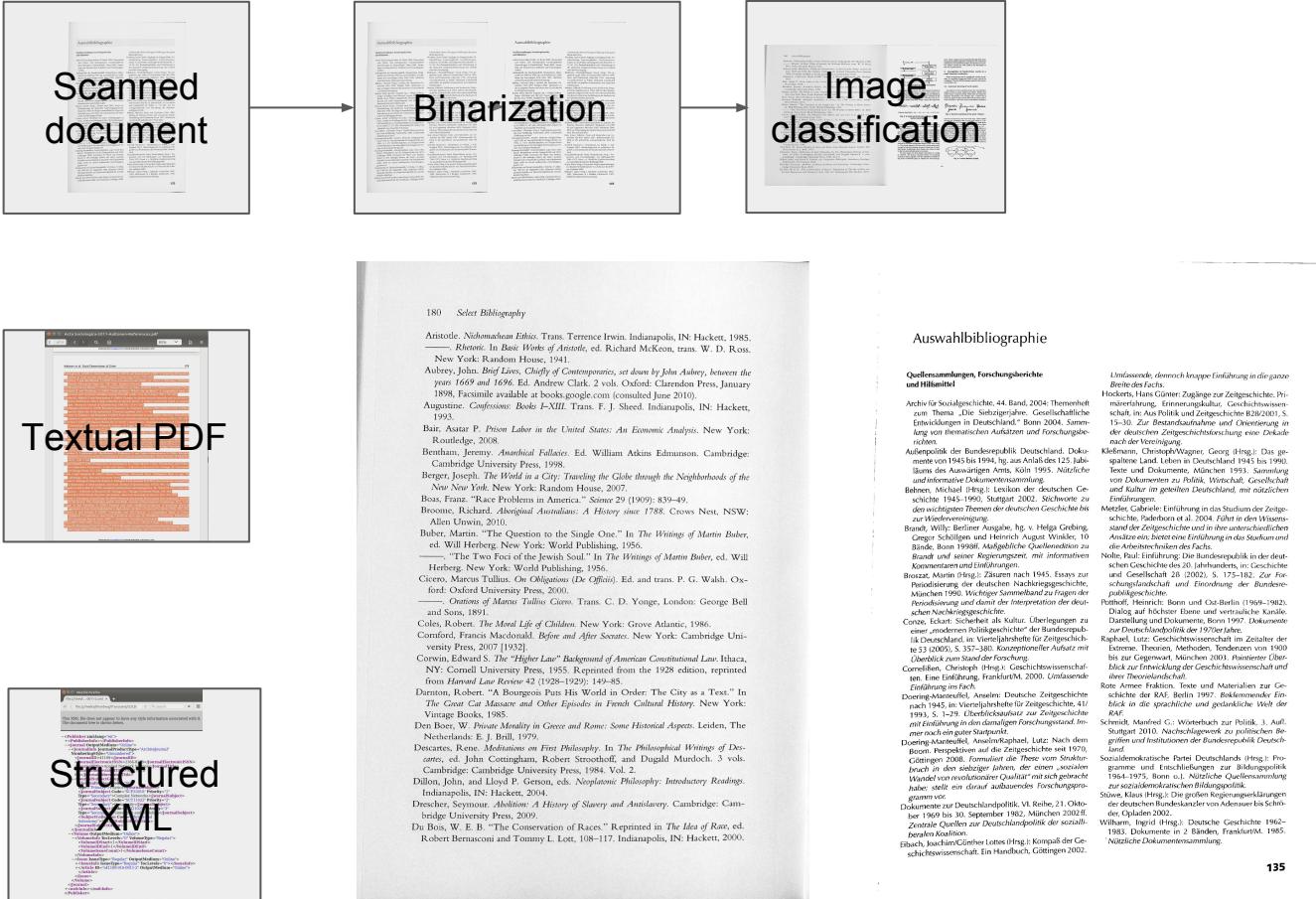


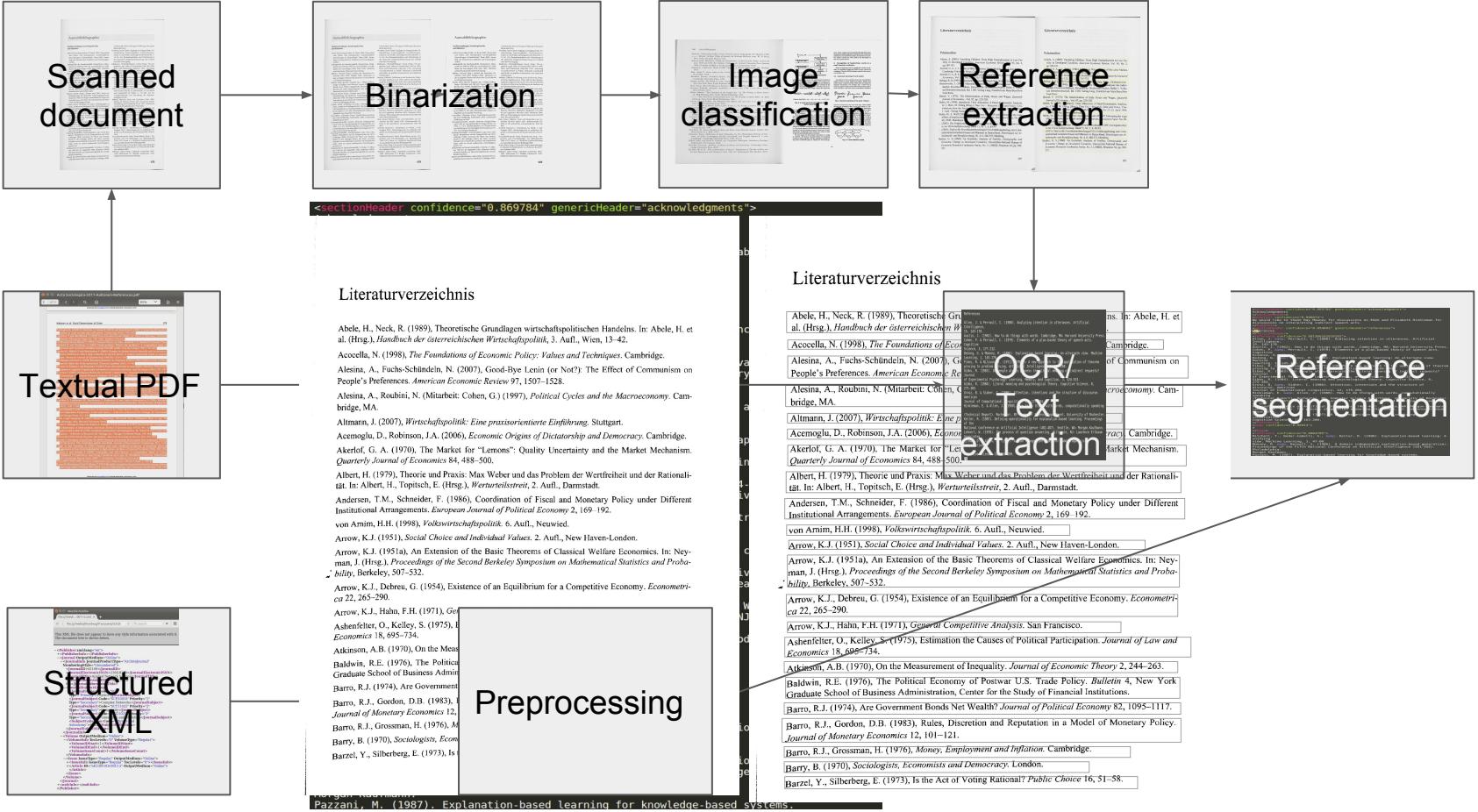
Structured
XML

This XML file does not appear to have any style information associated with it.
The document tree is shown below.

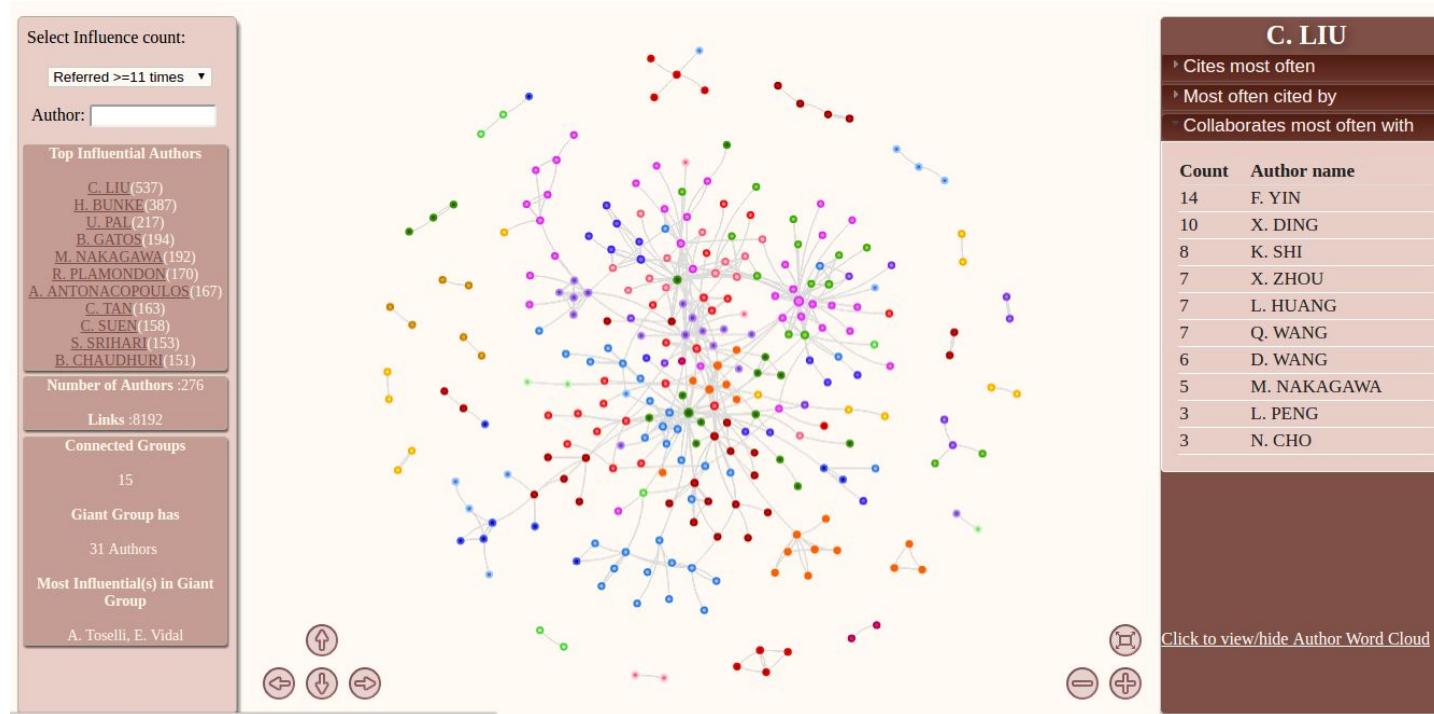
```
-<Publisher xml:lang="en">
+<PublisherInfo></PublisherInfo>
-<Journal OutputMedium="Online">
-<JournalInfo JournalProductType="ArchiveJournal"
NumberingStyle="Unnumbered">
<JournalID>41109</JournalID>
<JournalElectronicISSN>2364-8228</JournalElectronicISSN>
<JournalTitle>Applied Network Science</JournalTitle>
<JournalAbbreviatedTitle>Appl Netw
Sci</JournalAbbreviatedTitle>
-<JournalSubjectGroup>
<JournalSubject Code="SCP"
Type="Primary">Physics</JournalSubject>
<JournalSubject Code="SCP33010" Priority="1"
Type="Secondary">Complex Networks</JournalSubject>
<JournalSubject Code="SCT11022" Priority="2"
Type="Secondary">Complexity</JournalSubject>
<JournalSubject Code="SCI21025" Priority="3"
Type="Secondary">Simulation and Modeling</JournalSubject>
<SubjectCollection Code="Physics and
Astronomy">SC12</SubjectCollection>
</JournalSubjectGroup>
</JournalInfo>
-<Volume OutputMedium="Online">
-<VolumeInfo TocLevels="0" VolumeType="Regular">
<VolumeIDStart>1</VolumeIDStart>
<VolumeIDEnd>1</VolumeIDEnd>
<VolumeIssueCount>1</VolumeIssueCount>
</VolumeInfo>
-<Issue IssueType="Regular" OutputMedium="Online">
+<IssueInfo IssueType="Regular" TocLevels="0"></IssueInfo>
+<Article ID="s41109-016-0011-2" OutputMedium="Online">
</Article>
</Issue>
</Volume>
</Journal>
+<ns0:Info></ns0:Info>
</Publisher>
```



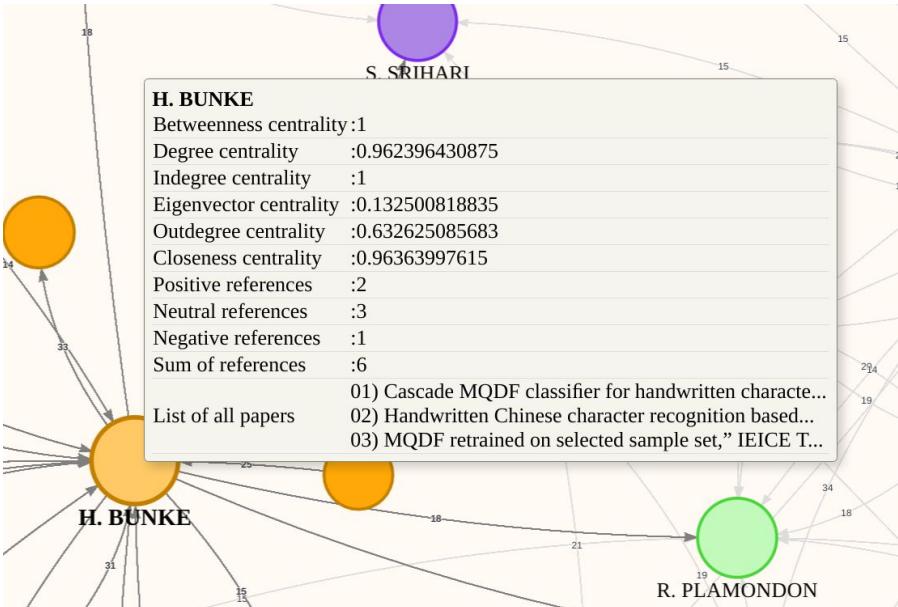
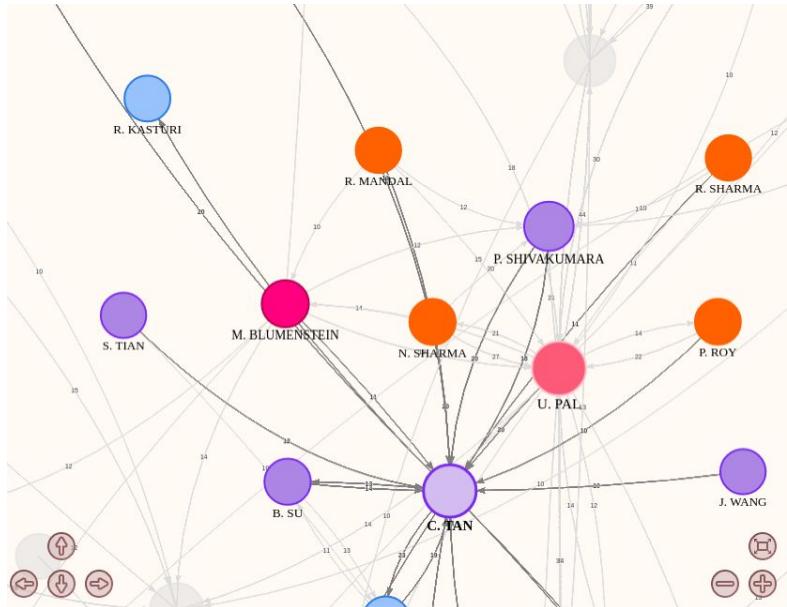




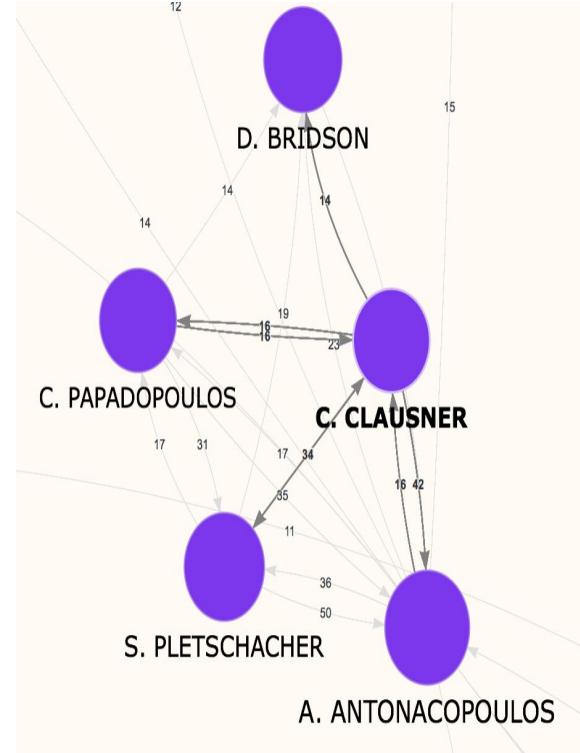
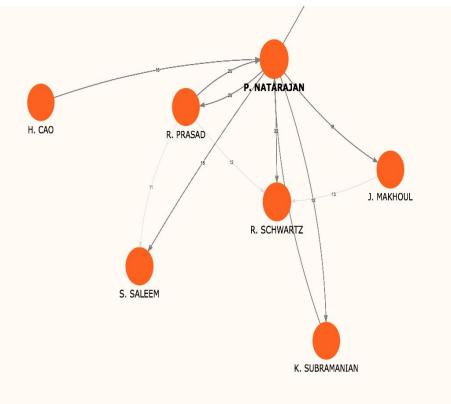
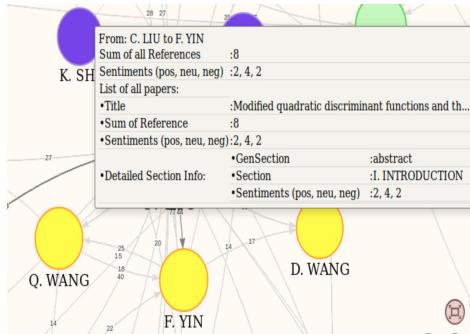
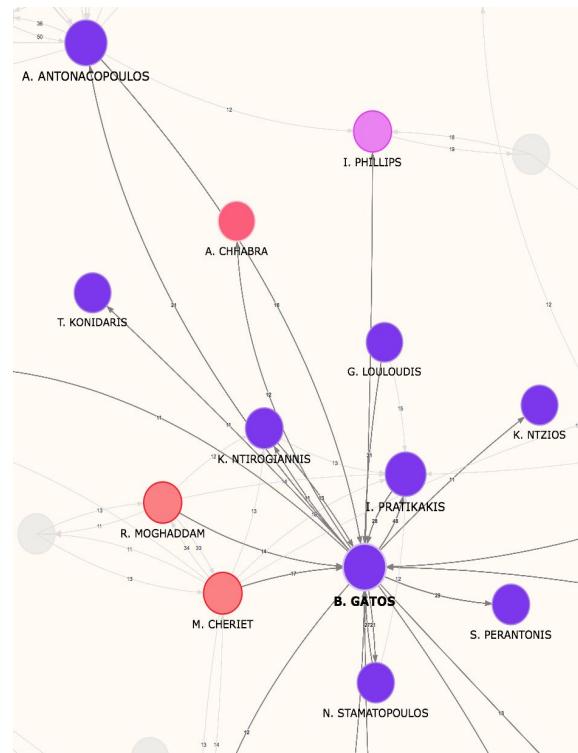
Appendix



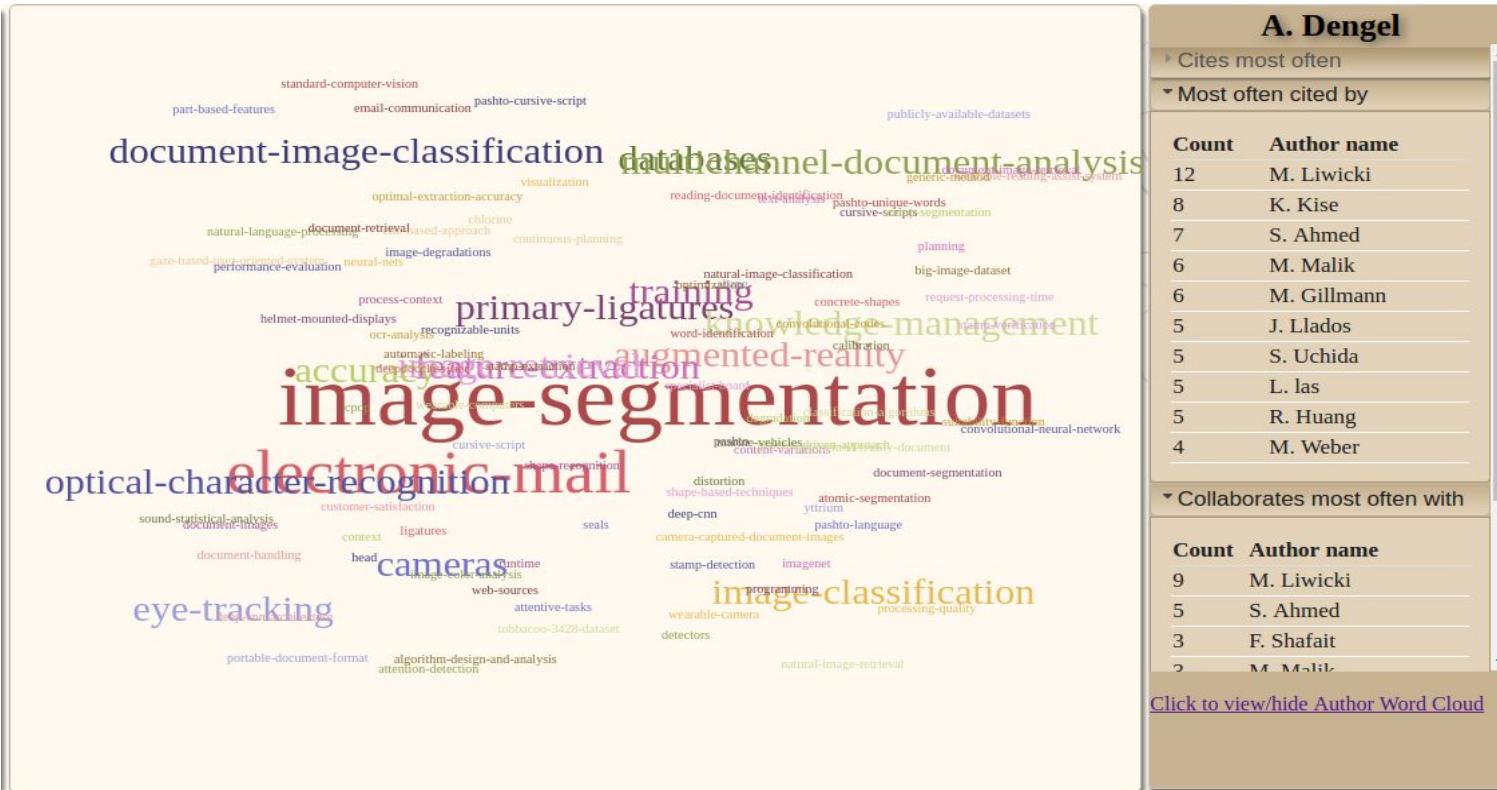
Appendix



Appendix



Appendix



Appendix

C. TAN

Cites most often

Count	Author name
67	C. TAN
23	S. LU
16	P. SHIVAKUMARA
14	A. JAIN
13	R. KASTURI
13	B. SU
11	I. PRATIKAKIS
11	B. GATOS
10	Y. ZHOU
9	H. BUNKE

Most often cited by

Count	Author name
67	C. TAN
20	N. SHARMA
20	P. SHIVAKUMARA
20	U. PAL
19	S. LU
14	B. SU
13	W. HUANG
12	M. CHERIET
12	S. TIAN
11	M. BLUMENSTEIN

Collaborates most often with

Count	Author name
11	S. LU
6	B. SU
6	P. SHIVAKUMARA
5	L. LI
4	S. TIAN
4	W. HUANG
4	Y. LU
3	B. YUAN
3	J. WANG
3	L. FAN

U. PAL

Cites most often

Count	Author name
109	U. PAL
44	B. CHAUDHURI
34	F. KIMURA
21	N. SHARMA
21	P. SHIVAKUMARA
20	C. TAN
20	M. BLUMENSTEIN
15	S. SRIHARI
14	R. PLAMONDON
14	P. ROY

Most often cited by

Count	Author name
109	U. PAL
31	B. CHAUDHURI
30	B. ROAD
30	F. KIMURA
27	N. SHARMA
22	P. ROY
19	M. BLUMENSTEIN
18	S. CHANDA
15	R. MANDAL
14	T. WAKABAYASHI

Collaborates most often with

Count	Author name
10	F. KIMURA
8	B. ROAD
7	B. CHAUDHURI
7	P. ROY
5	M. BLUMENSTEIN
5	N. SHARMA
4	A. ALEI
4	J. LLADOS
4	P. NAGABHUSHAN
4	T. WAKABAYASHI

V. GOVINDARAJU

Cites most often

Count	Author name
64	V. GOVINDARAJU
26	S. SRIHARI
16	Z. SHI
11	H. BUNKE
11	R. PLAMONDON
10	L. SCHOMAKER
9	G. KIM
9	S. SETLUR
9	C. RAMAIAH
8	A. JAIN

Most often cited by

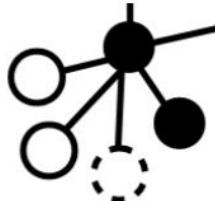
Count	Author name
64	V. GOVINDARAJU
20	C. RAMAIAH
20	S. SETLUR
20	Z. SHI
18	H. BUNKE
14	A. SHIVRAM
13	B. GATOS
12	P. NATARAJAN
11	S. SRIHARI
10	C. SUEN

Collaborates most often with

Count	Author name
10	S. SETLUR
8	Z. SHI
4	A. SHIVRAM
4	S. SRIHARI
3	C. RAMAIAH
3	S. MADHVANATH
3	S. TULYAKOV
2	B. ZHU
2	I. NWOGU
2	M. NAKAGAWA

Thank you.

<https://locdb.bib.uni-mannheim.de/blog/de/>



Tiessen, Jan (2007): Die Resultate im Blick?
ner/Döhler, Marian (Hrsg.): Agencies in W
Tondorf, Karin/Bahnmüller, Reinhard/Klaces,
instrument. Anwendungspraxis, Probleme
sigma.

Touraine, Alain (1984): Le retour de l'acteur: e
Treiber, Hubert (1984): Warum man nicht die
Mikroskop den ganzen Elefanten zu sehen.



Images

- [1] https://image.freepik.com/free-icon/male-user-shadow_318-34042.jpg
- [2] <https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcTuDltGNoXCvMPu3jwTMa-lio-DeK2TZljFZRK9THBduolliarvKA>
- [3] https://img.clipartfest.com/c123ecda18b92c3b0a147a994e57e0e8_to-do-list-clipart-things-to-do-list-clip-art_1024-1024.png
- [4] <https://scholar.google.de/>
- [5] <http://opencitations.net/>
- [6] https://us.intacct.com/sites/default/files/press-images/crossref_logo.png
- [7] <https://articles-images.sftcdn.net/wp-content/uploads/sites/8/2008/03/pdf.png>
- [8] <http://www.iconarchive.com/download/i22892/kyo-tux/phuzion/File-Web.ico>
- [9] https://encrypted-tbn2.gstatic.com/images?q=tbn:ANd9GcREEssBPV3h0H5JbJ1yrU_t5NRKhjyH5Hgnxd9iCkX3iExBnCXzAA
- [10] http://gsowww.gbv.de/images/logos/2_152.gif
- [11] <https://upload.wikimedia.org/wikipedia/en/thumb/e/eb/Springer.svg/1280px-Springer.svg.png>
- [12] <http://opencitations.net/static/img/logo.png>