

# “Challenges in Extracting and Managing References”

## EXCITE Workshop 2017

Philipp Mayr

GESIS – Leibniz Institute for the Social Sciences

2017-03-30



**#excitews2017**

<http://west.uni-koblenz.de/en/research/excite/workshop-2017>

## Who we are

- PI: Steffen Staab (WeST), Philipp Mayr (GESIS)
- Researchers: Behnam Ghavimi, Martin Körner
- Collaborator: Heinrich Hartmann (Independent)



## EXCITE: Background

- We run productive search systems and research in information retrieval, recommendation systems and knowledge discovery
  - Sowiport <http://sowiport.gesis.org/>
  - Related-work <http://dev.related-work.net/>
- **Shortage of citation data** for the international and German social sciences
- **Open availability of citation data is still very limited**
  - Open Citation Corpus <http://opencitations.net/>
  - CitEc: Citations in Economics <http://citec.repec.org/>

## EXCITE: Background

Service	Scope of the database	Coverage of the Social Sciences	Open Source/Data
Web of Science (WoS) <a href="http://webofknowledge.com">webofknowledge.com</a>	interdisciplinary	marginal	no/no
Scopus <a href="http://scopus.com">scopus.com</a>	interdisciplinary	marginal	no/no
Google Scholar (GS) <a href="http://scholar.google.com">scholar.google.com</a>	interdisciplinary	marginal	no/no
<a href="#">Microsoft Academic Search (MAS)</a> Deprecated in 2013	interdisciplinary	marginal	no/no
CiteSeerX <a href="http://citeseerx.ist.psu.edu">citeseerx.ist.psu.edu</a>	Computer Science, Mathematics, Physics	no coverage	yes/yes
ArnetMiner <a href="http://arnetminer.org">arnetminer.org</a>	Computer Science	no coverage	no/yes
CitEc <a href="http://citec.repec.org">citec.repec.org</a>	Economics	marginal	no/yes
Mendeley <a href="http://mendeley.com">mendeley.com</a>	interdisciplinary	no documentation	no/no
ResearchGate <a href="http://researchgate.net">researchgate.net</a>	interdisciplinary	no documentation	no/no
SSRN <a href="http://www.ssrn.com/">http://www.ssrn.com/</a>	Social Sciences	no documentation	no/no

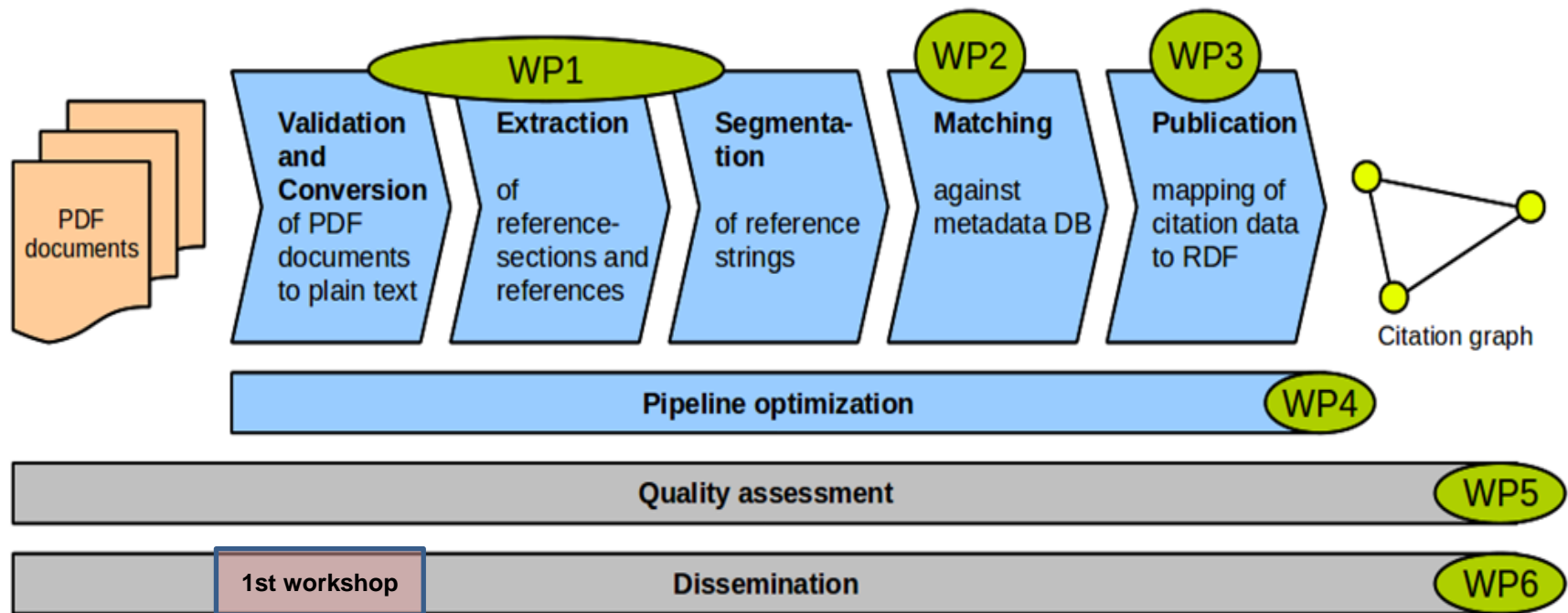
# EXCITE: Main objectives

- Narrow the supply gap of citation data in the social sciences
- Improve the breadth and accuracy of current reference extraction systems
- Develop web service API to allow third-parties to extract citation data from arbitrary publications
- Integrate and publish the extracted citation data

## General objectives

- Developing a toolchain of citation extraction and matching software
- Tools and data will be made available to researchers

## EXCITE: five steps



- (1) extraction of text from source documents,
- (2) identification of reference sections and other forms of embedded reference information within the text,
- (3) segmentation of individual references into its constituent fields such as author, title, etc.,
- (4) matching of reference strings against databases of bibliographic information,
- (5) the export of matched references to reusable formats

## EXCITE: Outcomes

- Open reference extraction tools for PDF documents
- Open datasets of segmented references
- Web service API for citation extraction
- Assessment of the overall quality of the extraction and matching pipeline
- Open gold standard testbed
- Improved infrastructures and services

## Workshop Agenda: Day 1

12:20	Information Extraction out of Born-Digital Scientific Articles	<a href="#">Roman Kern, TU Graz</a>
12:40	Advanced citation matching and large-scale full-text analysis	<a href="#">Nees Jan van Eck, Leiden U</a>
<b>13:00</b>	<b>Lunch Break (Cafeteria)</b>	
14:20	APIs for third parties to extract and deposit output executions of automated extraction pipelines (via <b>videoconferencing</b> )	<a href="#">Min-Yen Kan, NU Singapore</a>
14:40	Extracting references from scientific articles in CERMINE system	<a href="#">Dominika Tkaczyk, U Warsaw</a>
<b>15:00</b>	<b>Coffee Break (Cafeteria)</b>	
15:30	CitEc to CitEcCyr. A stab at distributed citation systems. (via <b>videoconferencing</b> )	<a href="#">Thomas Krichel, Open Library Society, NYC</a>
15:50	EXCITE project: Status report	<a href="#">Behnam Ghavimi, GESIS</a> <a href="#">Martin Körner, WeST</a> <a href="#">Heinrich Hartmann, Circonus</a>
16:10	Processing of in-text References: Towards a Semantic Analysis	<a href="#">Marc Bertin, U Toulouse</a>
16:30	Citations in Utopia Documents	<a href="#">David Thorne, U Manchester</a>
<b>16:50</b>	<b>Coffee Break (Cafeteria)</b>	
17:20	Research around the Tagging System BibSonomy	<a href="#">Andreas Hotho, U Würzburg</a>
17:50	LOC-DB: A Linked Open Citation Database provided by Libraries. Motivation and Challenges.	<a href="#">Kai Eckert, HDM Stuttgart</a> <a href="#">Anne Lauscher, HDM Stuttgart</a> <a href="#">Akansha Bhardwaj, DFKI</a> <sup>8</sup>
18:20	tbd (via <b>videoconferencing</b> )	<a href="#">Lee Giles, Penn State U</a>



## Workshop Agenda: Day 2

- Hand-on sessions

<b>9:00</b>	<b>Second Day Kickoff (Room: West II)</b>		
9:15	Extraction Result Discussion Group	Gold Standard Discussion Group	Collaboration Discussion Group
<b>11:15</b>	<b>Coffee Break (Cafeteria)</b>		
11:30	Extraction Result Discussion Group	Gold Standard Discussion Group	Collaboration Discussion Group
<b>12:30</b>	<b>Closing Talks (Room: West II)</b>		
<b>13:00</b>	<b>End</b>		

# Thank you

## Contact:

Dr Philipp Mayr

GESIS - Leibniz Institute for the Social Sciences, Germany

Email: [philipp.mayr@gesis.org](mailto:philipp.mayr@gesis.org)

Twitter: @philipp\_mayr

- Workshop website

<http://west.uni-koblenz.de/en/research/excite/workshop-2017>

- GIT <https://github.com/exciteproject/>

