

Joint Workshop on Bibliometric-enhanced Information Retrieval and Natural Language Processing for Digital Libraries (BIRNDL 2017)

Muthu Kumar Chandrasekaran

School of Computing,
National University of Singapore, Singapore
muthu.chandra@comp.nus.edu.sg

Kokil Jaidka

School of Arts & Sciences,
University of Pennsylvania, USA
jaidka@sas.upenn.edu

Philipp Mayr

GESIS – Leibniz Institute for the Social
Sciences, Germany
philipp.mayr@gesis.org

ABSTRACT

The large scale of scholarly publications poses a challenge for scholars in information seeking and sensemaking. Bibliometrics, information retrieval (IR), text mining and NLP techniques could help in these search and look-up activities, but are not yet widely used. This workshop is intended to stimulate IR researchers and digital library professionals to elaborate on new approaches in natural language processing, information retrieval, scientometrics, text mining and recommendation techniques that can advance the state-of-the-art in scholarly document understanding, analysis, and retrieval at scale. The BIRNDL workshop at SIGIR 2017 will incorporate an invited talk, paper sessions and the third edition of the Computational Linguistics (CL) Scientific Summarization Shared Task.

CCS CONCEPTS

• **Information systems** → **Information retrieval**; *Link and co-citation analysis*; • **Applied computing** → **Digital libraries and archives**;

KEYWORDS

Scientometrics; Information Retrieval; Digital Libraries; NLP; Summarization; Information Extraction; Citation analysis

ACM Reference format:

Muthu Kumar Chandrasekaran, Kokil Jaidka, and Philipp Mayr. 2017. Joint Workshop on Bibliometric-enhanced Information Retrieval and Natural Language Processing for Digital Libraries (BIRNDL 2017). In *Proceedings of SIGIR'17, August 7–11, 2017, Shinjuku, Tokyo, Japan.*, 2 pages. <https://doi.org/10.1145/3077136.3084370>

1 INTRODUCTION

Over the past several years, the BIRNDL workshop and its parent workshops are establishing themselves as the primary interdisciplinary venue for the cross-pollination of bibliometrics and information retrieval (IR) [5]. Our motivation as organizers of the workshop started from the observation that both communities share only a partial overlap; yet, the main discourse in both fields consists of different approaches to solve similar problems. We believe

that a knowledge transfer would be profitable for both sides. A good overview of the symbiotic relationship that exists among bibliometrics, IR and natural language processing (NLP) has been presented by Wolfram [6]. A report of the past BIRNDL workshop has been published recently in The SIGIR Forum [1].

The goal of the BIRNDL workshop at SIGIR is to engage the IR community about the open problems in academic search. Academic search refers to the large, cross-domain digital repositories which index research papers, such as the ACL Anthology, ArXiv, ACM Digital Library, IEEE database, Web of Science and Google Scholar. Currently, digital libraries collect and allow access to papers and their metadata — including citations — but mostly do not analyze the items they index. The scale of scholarly publications poses a challenge for scholars in their search for relevant literature. Finding relevant scholarly literature is the key theme of BIRNDL and sets the agenda for tools and approaches to be discussed and evaluated at the workshop.

Papers at the 2nd BIRNDL workshop will incorporate insights from IR, bibliometrics and NLP to develop new techniques to address the open problems such as evidence-based searching, measurement of research quality, relevance and impact, the emergence and decline of research problems, identification of scholarly relationships and influences and applied problems such as language translation, question-answering and summarization. We will also address the need for established, standardized baselines, evaluation metrics and test collections. Towards the purpose of evaluating tools and technologies developed for digital libraries, we are organizing the 3rd CL-SciSumm Shared Task based on the CL-SciSumm corpus, which comprises over 500 computational linguistics (CL) research papers, interlinked through a citation network.

The organizers of the 2nd BIRNDL workshop at SIGIR 2017¹ have previously organized other workshop series at premier IR and CS venues - notably, the Bibliometric-enhanced Information Retrieval (BIR) workshops in 2014, 2015 and 2016 at ECIR [4] and the NLP4DL workshop at ACL-IJCNLP (2009). Most recently, the BIRNDL workshop and the 2nd CL-SciSumm Shared Task were co-located with JCDL 2016² [1], where 10 research papers and 10 system papers were presented³ (acceptance rate: 30%). In 2017, the BIRNDL workshop takes this legacy forward with a focus on scholarly publications and data, and an updated scientific summarization Shared Task for its participants.

This workshop will be relevant to scholars in computer and information science, specializing in IR and NLP. It will also be of

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

SIGIR'17, August 7–11, 2017, Shinjuku, Tokyo, Japan.

© 2017 Copyright held by the owner/author(s).

ACM ISBN ACM ISBN 978-1-4503-5022-8/17/08.

<https://doi.org/10.1145/3077136.3084370>

¹<http://wing.comp.nus.edu.sg/birndl-sigir2017/>

²<http://wing.comp.nus.edu.sg/birndl-jcdl2016/>

³<http://ceur-ws.org/Vol-1610/>

importance to all stakeholders in the publication pipeline: practitioners, publishers and policymakers. Today’s publishers continue to provide new ways to support their consumers in disseminating and retrieving the right published works to their audience. Formal citation metrics are increasingly a factor in decision-making by universities and funding bodies worldwide, making the need for research in applying these metrics more pressing.

2 WORKSHOP TOPICS AND FORMAT

Our goal is to encourage insights from IR, NLP and CL for scholarly document understanding, document analysis and retrieval in digital libraries. The papers presented at the workshop will touch upon several topics, including (but not limited to) full-text analysis, multimedia and multilingual analysis and alignment as well as the application of citation-based NLP, information retrieval and information seeking techniques in digital libraries. More specifically, our fields of interests include:

- Infrastructures for scientific text mining and IR
- Semantic and Network-based indexing, navigation, searching and browsing in structured data
- Discourse structure identification and argument mining from scientific papers
- Summarization and question-answering for scholarly DLs
- Recommendation for scholarly papers, reviewers, citations and publication venues
- Measurement and evaluation of quality and impact
- Metadata and controlled vocabularies for resource description and discovery; automatic metadata discovery, such as language identification
- Disambiguation issues in scholarly DLs using NLP or IR techniques; data cleaning and data quality.

2.1 Tentative Schedule of Events

The workshop will start with a keynote titled “Do “Future Work” sections have a real purpose? Citation links and entailment for global scientometric questions” by Dr. Simone Teufel (University of Cambridge). This session will be followed by regular research paper presentations, overview papers and posters on the Shared Task.

2.2 The CL-SciSumm Shared Task

The 3rd Computational Linguistics (CL) Scientific Summarization Shared Task, sponsored by Microsoft Research Asia, will be conducted as a part of this workshop. This is the first medium-scale shared task on scientific document summarization in the CL domain. It follows up on and extends the successful CL Shared Tasks conducted as a part of BIRNDL 2016 [1], and within the Biomed-Summ Track at the Text Analysis Conference 2014 (TAC 2014) [2]. In the CL-SciSumm 2016 [3] Shared Task, fifteen teams from six countries signed up, and ten teams ultimately submitted and presented their results.

The Shared Task comprises three sub-tasks in automatic research paper summarization on a new corpus of research papers, as described below.

Given: A topic consisting of a Reference Paper (RP) and up to ten Citing Papers (CPs) that all contain citations to the RP. Citations in

the CP are pre-identified as the text spans (i.e., citances), that cite the RP.

Task 1a: For each citance, identify the spans of text (cited text spans) in the RP that most accurately reflect the citance.

Task 1b: For each cited text span, identify what facet of the paper it belongs to, from a predefined set of facets.

Task 2 (optional bonus task): Finally, generate a structured summary of the RP from the cited text spans of the RP. The length of the summary should not exceed 250 words.

Evaluation: Task 1 will be scored by overlap of text spans measured by number of sentences in the system output vs gold standard. Task 2 will be scored using the ROUGE family of metrics between the system output, and i) human summaries, ii) community summaries comprising the cited text spans, and ii) the Abstract section of the reference paper.

This task is continues to be of interest to a broad community including those working in CL and NLP, especially in the sub-disciplines of text summarization, discourse structure in scholarly discourse, paraphrase, textual entailment and text simplification.

3 OUTLOOK

This workshop is the first step to foster a reflection on interdisciplinarity, and the benefits that the disciplines Bibliometrics, IR and NLP can derive from it in the Digital Libraries context. The authors of accepted papers will be invited to submit extended versions of their work to the International Journal on Digital Libraries (IJDL). As an output of BIRNDL 2016, a special issue of IJDL on “Bibliometrics, Information Retrieval and Natural Language Processing in Digital Libraries” is currently in preparation. In the future, we plan to continue to host this series of workshops and Shared Tasks at prominent IR, NLP and Digital Library venues.

ACKNOWLEDGMENTS

We thank Microsoft Research Asia for their generous support in funding the development, dissemination and organization of the CL-SciSumm dataset and the Shared Task. We are also grateful to the co-organizers of the 1st BIRNDL workshop - Guillaume Cabanac, Ingo Frommholz, Min-Yen Kan and Dietmar Wolfram, for their continued support and involvement.

REFERENCES

- [1] Guillaume Cabanac, Muthu Kumar Chandrasekaran, Ingo Frommholz, Kokil Jaidka, Min-Yen Kan, Philipp Mayr, and Dietmar Wolfram. 2016. Report on the Joint Workshop on Bibliometric-enhanced Information Retrieval and Natural Language Processing for Digital Libraries (BIRNDL 2016). *SIGIR Forum* 50, 2 (2016), 36–43. <http://sigir.org/wp-content/uploads/2017/01/p036.pdf>
- [2] Kokil Jaidka, Muthu Kumar Chandrasekaran, Beatriz Fisas Elizalde, Rahul Jha, Christopher Jones, Min-Yen Kan, Ankur Khanna, Diego Molla-Alíod, Dragomir Radev, Francesco Ronzano, et al. 2014. The computational linguistics summarization pilot task. In *Proceedings of Text Analysis Conference*. Gaithersburg, USA.
- [3] Kokil Jaidka, Muthu Kumar Chandrasekaran, Sajal Rustagi, and Min-Yen Kan. 2016. Overview of the CL-SciSumm 2016 Shared Task. In *In Proceedings of Joint Workshop on Bibliometric-enhanced Information Retrieval and NLP for Digital Libraries (BIRNDL 2016)*.
- [4] Philipp Mayr, Ingo Frommholz, and Guillaume Cabanac. 2016. Report on the 3rd International Workshop on Bibliometric-enhanced Information Retrieval (BIR 2016). *SIGIR Forum* 50, 1 (2016), 28–34. <http://sigir.org/files/forum/2016/p028.pdf>
- [5] Philipp Mayr and Andrea Scharnhorst. 2015. Scientometrics and Information Retrieval - weak-links revitalized. *Scientometrics* 102, 3 (2015), 2193–2199. <https://doi.org/10.1007/s11192-014-1484-3>
- [6] Dietmar Wolfram. 2016. Bibliometrics, Information Retrieval and Natural Language Processing: Natural Synergies to Support Digital Library Research. In *Proc. of the BIRNDL Workshop 2016*. 6–13. <http://ceur-ws.org/Vol-1610/paper1.pdf>