

**Proceedings of the
Computational Sanskrit & Digital Humanities**

Selected papers presented
at the 18th World Sanskrit Conference

(in online mode)
Canberra, Australia
January 9 – 13, 2023

Edited by
AMBA KULKARNI & OLIVER HELLWIG

©2022 The Association for Computational Linguistics
Order copies of this and other ACL proceedings from:
Association for Computational Linguistics (ACL) 209 N. Eighth Street
Stroudsburg, PA 18360
USA
Tel: +1-570-476-8006 Fax: +1-570-476-0860 acl@aclweb.org

ISBN 978-1-950737-71-0

Preface

This volume presents edited versions of shortlisted papers accepted for presentation in the session ‘Computational Sanskrit and Digital Humanities’ at the 18th World Sanskrit Conference during Jan 9-13, 2023 in online mode. The physical conference was delayed by two years due to the global pandemic situation. Thus there were two rounds of calls for submissions. In order to disseminate the research among the scholars, the convenors of the Computational Sanskrit and Digital Humanities decided to hold an online event after the first round of submission. We thank the organisers of the World Sanskrit Conference and in particular the International Association of Sanskrit Studies for giving us permission to hold this event online. The papers shortlisted in the first round were presented at the online event during Jan 11-12, 2022, as well.

We received a total of 29 submissions in the two rounds. Each submission was reviewed by at least three reviewers and 15 submissions were shortlisted for presentation.

A good number of papers focus on how machine learning approaches can be applied to Sanskrit texts, paying special attention to its status as a low resource language. The contribution of Krishna et al. is an in-depth study of data-driven parsing algorithms for Classical Sanskrit. It discusses various parsing architectures, the choice of input features and possible sources of parsing errors. Another aspect of data-driven NLP is covered in the contribution of Sandhan et al. who compare static and contextualized word embeddings on a range of semantic classification tasks and obtain remarkably good results given the small size of the digital Sanskrit corpus. As a foundation for data-driven methods, Krishnan et al. discuss endeavors towards merging existing annotated corpora of Classical Sanskrit in a common format. Building NLP resources is also the topic of the paper by Sarkar et al. who develop a simple, yet effective machine learning approach for pre-annotating an NER data set. While the papers mentioned so far use quantitative models for language analysis, Maity et al. discuss the problem of disambiguating the kāraka notion of oblique cases from a Pāṇinian perspective and evaluate their approach on a small sample of annotated sentences.

The remaining papers of this volume cover aspects of (Sanskrit) NLP that go beyond the identification and analysis of basic linguistic structures. Mahesh and Bhattacharya evaluate how well high-level textual features can be predicted with contextualized word embeddings. The contribution of Neill describes the open-source project Skrutable, a tool for identifying meters in Sanskrit texts, and sketches directions for future research in this field. Another paper focusing on meter identification is presented by Terdalkar and Bhattacharya whose system also includes a fuzzy search option. Ajotikar et al. propose an extension of the TEI standard for encoding various levels of information present in the Sanskrit commentarial literature and discuss which research questions can be addressed by applying XSLT templates to textual sources encoded in this way. Another contribution by Scharf et al. describes the TEI encoding of the Rāmopakhyāna and the Kramapha system used to interactively display its content in the web. Terdalkar et al. report about ongoing work on building knowledge graphs from Āyurvedic texts. Their paper elaborates on the annotation process and its prerequisites as well as on the design of an appropriate medical and botanical ontology. The problem of extracting and representing the knowledge structure of śāstra texts is also addressed in the paper by Susarla et al. who concentrate on sentence level features of scientific texts. Hellwig et al. describe an extension of Bloomfield’s concordance of Vedic mantras and discuss its application to open issues in Vedic Studies.

Finally, two papers address aspects of Middle Indo-Aryan languages and the Buddhist literature composed in them. Harnsukworapanich et al. give an overview of D-Tipitaka, an online version of

the Dhammachai Tipiṭaka Edition with linked manuscript images and contextual information. Zigmund discusses how Pāli commentaries can be distinguished from the commented canonical texts by applying standard clustering algorithms to frequent words in these texts.

We thank the Convenors, Programme Committee members and the numerous experts who helped us in the review process, and all our authors who responded positively to the reviewer's comments and improved their manuscripts accordingly. We thank the entire 18th WSC organising committee, led by Prof McComas Taylor, which provided us the necessary logistic support for the organisation of this section.

Amba Kulkarni & Oliver Hellwig

Programme Committee

- **Convenors:**

- Gérard Huet (Inria, Paris, France) (Chair)
- Amba Kulkarni (University of Hyderabad, India)

- **Chairs:**

- Amba Kulkarni (University of Hyderabad, India)
- Oliver Hellwig (University of Zurich, Germany)

- **Members:**

- Ivan Andrijani (University of Zagreb, Croatia)
- Stefan Baums (University of Munich, Germany)
- Arnab Bhattacharya (IIT Kanpur, India)
- Brendan Gillon (McGill University, Canada)
- Pawan Goyal (IIT Kharagpur, India)
- Malhar Kulkarni (IIT Bombay, India)
- Dhaval Patel (Ahmedabad, India)
- Wiebke Petersen (University of Düsseldorf, Germany)
- Pavan Kumar Satuluri (Chinmaya Vishwavidyapeeth, Veliyanad, India)
- Sai Susarla (MIT, Pune, India)
- Peter Scharf (IIIT Hyderabad, India)
- Srinivasa Varakhedi (KKSU, Ramtek, India)

Table of Contents

<i>Neural Approaches for Data Driven Dependency Parsing in Sanskrit</i> Amrith Krishna, Ashim Gupta, Deepak Garasangi, Jivnesh Sandhan, Pavankumar Satuluri and Pawan Goyal	1
<i>Evaluating Neural Word Embeddings for Sanskrit</i> Jivnesh Sandhan, Om Adideva Paranjay, Digumarthi Komal, Laxmidhar Behera and Pawan Goyal ...	21
<i>Validation and Normalization of DCS corpus and Development of the Sanskrit Heritage Engine's Segmenter</i> Sriram Krishnan, Amba Kulkarni and Gérard Huet	38
<i>Pre-annotation Based Approach for Development of a Sanskrit Named Entity Recognition Dataset</i> Sujoy Sarkar, Amrith Krishna and Pawan Goyal	59
<i>Disambiguation of Instrumental, Dative and Ablative Case suffixes in Sanskrit</i> Malay Maity, Sanjeev Panchal and Amba Kulkarni	71
<i>Creation of a Digital Rig Vedic Index (Anukramani) for Computational Linguistic Tasks</i> A V S D S Mahesh and Arnab Bhattacharya	89
<i>Skrutable: Another Step Toward Effective Sanskrit Meter Identification</i> Tyler Neill	97
<i>Chandojnanam: A Sanskrit Meter Identification and Utilization System</i> Hrishikesh Terdalkar and Arnab Bhattacharya	113
<i>Development of a TEI standard for digital Sanskrit texts containing commentaries: A pilot study of Bhaṭṭi's Rāvaṇavadha with Mallinātha's commentary on the first canto</i> Tanuja P A jotikar and Peter M Scharf	128
<i>Rāmopākhyāna: A Web-based reader and index</i> Peter M Scharf and Dhruv Chauhan	146
<i>Semantic Annotation and Querying Framework based on Semi-structured Ayurvedic Text</i> Hrishikesh Terdalkar, Arnab Bhattacharya, Madhulika Dubey, Ramamurthy S and Bhavna Naneria Singh	155
<i>Shaastra Maps: Enabling Conceptual Exploration of Indic Shaastra Texts</i> Sai Susarla, Suryanarayana Jammalamadaka, Vaishnavi Nishankar, Siva Panuganti, Anupama Ryali and Sushrutha S	174
<i>The Vedic corpus as a graph. An updated version of Bloomfields Vedic Concordance</i> Oliver Hellwig, Sven Sellmer and Kyoko Amano	188

<i>The transmission of the Buddha's teachings in the digital age</i>	
Sumachaya Harnsukworapanich and Phatchareporn Suphipat	201
<i>Distinguishing Commentary from Canon: Experiments in Pāli Computational Linguistics</i>	
Dan Zigmund	213