

RANLPStud 2021

**Proceedings of the
Student Research Workshop**

associated with

**The 13th International Conference on
Recent Advances in Natural Language Processing
RANLP'2021**

1–3 September, 2021
(held online)

STUDENT RESEARCH WORKSHOP
ASSOCIATED WITH THE INTERNATIONAL CONFERENCE
RECENT ADVANCES IN
NATURAL LANGUAGE PROCESSING'2021

PROCEEDINGS

1-3 September 2021

Series Online ISSN 2603-2821

Designed by INCOMA Ltd.
Shoumen, BULGARIA

Preface

The RANLP 2021 Student Research Workshop (RANLPStud 2021) is a special track of the established international conference Recent Advances in Natural Language Processing (RANLP 2021), now in its thirteenth edition as an online event due to the unusual circumstances related to the global pandemic.

The RANLP Student Research Workshop is being organised for the seventh time and this year is running in parallel with the other tracks of the main RANLP 2021 conference. The target of RANLPStud 2021 is to be a discussion forum and provide an outstanding opportunity for students at all levels (Bachelor, Masters, and Ph.D.) to present their work in progress or completed projects to an international research audience and receive feedback from senior researchers.

The RANLP 2021 Student Research Workshop received a large number of submissions, this year thirty-three (33) papers were submitted to the event covering nearly all the continents: Asia, The Americas (North and South), Europe and Oceania, a fact which was reflecting the record number of events, sponsors, submissions, and participants at the main RANLP 2021 conference.

We have accepted four (4) excellent student papers as long presentations, one of them has received the Best Paper Award in the opening session also attended by mentors from academia and 29 submissions are presented in a shorter form.

We did our best to make the reviewing process in the interest of our authors, by asking our reviewers to give as exhaustive comments and suggestions as possible as well as to maintain an encouraging attitude. Each student submission was reviewed by at least 2 Programme Committee members, which are specialists in their field and were carefully selected to match the submission's topic.

This year, as usual, we invited both strictly Natural Language Processing (NLP) papers, and submissions at the borderline between two sciences (but bearing contributions to NLP. The topics of the accepted submissions include: Computational treatment of humour, sarcasm and irony; Computer-aided language learning; Computational cognitive modelling ; Corpus annotation; Crowdsourcing; Deep learning for NLP; Discourse text summarisation; Fact checking; Information extraction; Language resources and corpora; Lexicography; Lexicon; Machine translation, including statistical machine translation and neural machine translation; NLP for biomedical texts; NLP for social media; Natural language generation; Natural language processing for educational applications; Opinion mining and sentiment analysis; Pragmatics; Phonetics and phonology; Semantics; Similarity; Syntax ; Text and web mining; Translation technology including translation memory systems; Word embeddings; Word-sense disambiguation.

We are thankful to the members of the Programme Committee for having provided such exhaustive reviews and even accepting additional reviews, and to the conference mentors, who provided additional comments to participants.

We would like especially to thank all members of the Organising Committee (listed in alphabetical order) for their long and devoted work : Souhila Djabri, Dinara Gimadi, Tsvetomila Mihaylova and Ivelina Nikolova-Koleva. THANK YOU for making this event possible within a very short time limit.

The RANLPStud 2021 Organisers

Souhila Djabri, University of Alicante, Spain

Dinara Gimadi, University of Wolverhampton, United Kingdom

Tsvetomila Mihaylova, Institute of Telecommunications, Lisbon, Portugal

Ivelina Nikolova-Koleva, Bulgarian Academy of Sciences and Sirma AI, Bulgaria

Organizers:

Souhila Djabri (University of Alicante, Spain)
Dinara Gimadi (University of Wolverhampton, United Kingdom)
Tsvetomila Mihaylova (Institute of Telecommunications, Lisbon, Portugal)
Ivelina Nikolova-Koleva (Bulgarian Academy of Sciences and Sirma AI, Bulgaria)

Programme Committee:

Aida Kostikova (New Bulgarian University, Bulgaria)
Ali Hatami (University of Wolverhampton, United Kingdom)
Alistair Plum (University of Wolverhampton, United Kingdom)
Andrey Tagarev (Sirma AI, Bulgaria)
Burcu Can Buglalilar (University of Wolverhampton, United Kingdom)
Cengiz Acartürk (Middle East Technical University, Turkey)
Constantin Orăsan (University of Surrey, United Kingdom)
Darya Filippova (University of Wolverhampton, United Kingdom)
Dean Hunter (University of Wolverhampton, United Kingdom)
Dinara Gimadi (New Bulgarian University, Bulgaria)
Elena Murgolo (Aglatech14, Italy)
Emma Franklin (University of Wolverhampton, United Kingdom)
Frédéric Blain (University of Wolverhampton, United Kingdom)
Georgi Karadzhov (University of Cambridge, United Kingdom)
Georgi Georgiev (Releva.ai, Bulgaria)
Ivan Koychev (Sofia University, Bulgaria)
Ivelina Nikolova-Koleva (Bulgarian Academy of Sciences and Sirma AI, Bulgaria)
Jessica López Espejel (Université Sorbonne Paris Nord, France)
Le An Ha (University of Wolverhampton, United Kingdom)
Maria Kunilovskaya (Research Group in Computational Linguistics, United Kingdom)
Maria Carmela Cariello (University of Pisa, Italy)
Marie Escribe (University of Wolverhampton, United Kingdom)
Momchil Hardalov (Sofia Univeristy, Bulgaria)
Necva Bölücü (Hacettepe University, Turkey)
Nikola Spasovski (University of Wolverhampton, United Kingdom)
Preslav Nakov (Qatar Computing Research Institute, HBKU, Qatar)
Raheem Sarwar (University of Wolverhampton, United Kingdom)
Richard Evans (University of Wolverhampton, United Kingdom)
Sandra Kübler (Indiana University, United States)
Sara Može (University of Wolverhampton, United Kingdom)
Shaifali Khulbe (University of Wolverhampton, United kingdom)
Sonia Kropiowska (University of Wolverhampton, United Kingdom)
Souhila Djabri (University of Alicante, Spain)
Svetla Boytcheva (Bulgarian Academy of Sciences and Sirma AI, Bulgaria)
Svetla Koeva (Bulgarian Academy of Sciences, Bulgaria)
Tharindu Ranasinghe (University of Wolverhampton, United Kingdom)
Tsvetomila Mihaylova (Institute of Telecommunications, Lisbon, Portugal)
Yasen Kiproff (Sofia University, Bulgaria)

Table of Contents

<i>Humor Generation and Detection in Code-Mixed Hindi-English</i> Kaustubh Agarwal and Rhythm Narula	1
<i>Towards Code-Mixed Hinglish Dialogue Generation</i> Vibhav Agarwal, Pooja Rao and Dinesh Babu Jayagopi	7
<i>Hinglish to English Machine Translation using Multilingual Transformers</i> Vibhav Agarwal, Pooja Rao and Dinesh Babu Jayagopi	16
<i>SiPOS: A Benchmark Dataset for Sindhi Part-of-Speech Tagging</i> Wazir Ali, Zenglin Xu and Jay Kumar	22
<i>Sarcasm Detection and Building an English Language Corpus in Real Time</i> Oliver Cakebread-Andrews	31
<i>Correcting Texts Generated by Transformers using Discourse Features and Web Mining</i> Alexander Chernyavskiy, Dmitry Ilvovsky and Boris Galitsky	36
<i>Introducing linguistic transformation to improve translation memory retrieval. Results of a professional translators' survey for Spanish, French and Arabic</i> Souhila Djabri and Rocío Caro Quintana	44
<i>Using Transfer Learning to Automatically Mark L2 Writing Texts</i> Tim Elks	51
<i>Bilingual Terminology Extraction Using Neural Word Embeddings on Comparable Corpora</i> Darya Filippova, Burcu Can and Gloria Corpas Pastor	58
<i>Web-sentiment analysis of public comments (public reviews) for languages with limited resources such as the Kazakh language</i> Dinara Gimadi	65
<i>Disambiguating Grammatical Number and Gender With BERT</i> Annegret Janzso	69
<i>Towards a Language Model for Temporal Commonsense Reasoning</i> Mayuko Kimura, Lis Kanashiro Pereira and Ichiro Kobayashi	78
<i>Text Preprocessing and its Implications in a Digital Humanities Project</i> Maria Kunilovskaya and Alistair Plum	85
<i>Compiling a specialised corpus for translation research in the environmental domain</i> Anastasiia Laktionova	94
<i>Paragraph Similarity Matches for Generating Multiple-choice Test Items</i> Halyna Maslak and Ruslan Mitkov	99
<i>Neural Borrowing Detection with Monolingual Lexical Models</i> John Miller, Emanuel Pariasca and Cesar Beltran Castañón	109
<i>Does local pruning offer task-specific models to learn effectively ?</i> Abhishek Kumar Mishra and Mohna Chakraborty	118

<i>On Reducing Repetition in Abstractive Summarization</i>	
Pranav Nair and Anil Kumar Singh	126
<i>Improving Abstractive Summarization with Commonsense Knowledge</i>	
Pranav Nair and Anil Kumar Singh	135
<i>A Dataset for Research on Modelling Depression Severity in Online Forum Data</i>	
Isuri Anuradha Nanomi Arachchige, Vihangi Himaya Jayasuriya and Ruvan Weerasinghe	144
<i>Handling synset overgeneration: Sense Merging in BTB-WN</i>	
Ivaylo Radev and Zara Kancheva	154
<i>On the Evolution of Word Order</i>	
Idan Rejwan and Avi Caciularu	162
<i>EmoPars: A Collection of 30K Emotion-Annotated Persian Social Media Texts</i>	
Nazanin Sabri, Reyhane Akhavan and Behnam Bahrak	167
<i>A Review on Document Information Extraction Approaches</i>	
Kanishka Silva and Thushari Silva	174
<i>Towards New Generation Translation Memory Systems</i>	
Nikola Spasovski and Ruslan Mitkov	180
<i>Question answering in Natural Language: the Special Case of Temporal Expressions</i>	
Armand Stricker	184
<i>Building A Corporate Corpus For Threads Constitution</i>	
Lionel Tadonfouet Tadjou, Fabrice Bourge, Tiphaine Marie, Laurent Romary and Éric de la Clergerie	193
<i>Generating Answer Candidates for Quizzes and Answer-Aware Question Generators</i>	
Kristiyan Vachev, Momchil Hardalov, Georgi Karadzhov, Georgi Georgiev, Ivan Koychev and Preslav Nakov	203
<i>Toward Discourse-Aware Models for Multilingual Fake News Detection</i>	
Francielle Vargas, Fabrício Benevenuto and Thiago Pardo	210
<i>Automatic Transformation of Clinical Narratives into Structured Format</i>	
Sylvia Vassileva, Gergana Todorova, Kristina Ivanova, Boris Velichkov, Ivan Koychev, Galia Angelova and Svetla Boytcheva	219