

Biomedical NLP Workshop 2017

**Proceedings of the  
Biomedical NLP Workshop**

*associated with*

**The 11th International Conference on  
Recent Advances in Natural Language Processing  
(RANLP 2017)**

8 September, 2017  
Varna, Bulgaria

BIOMEDICAL NATURAL LANGUAGE PROCESSING WORKSHOP  
ASSOCIATED WITH THE INTERNATIONAL CONFERENCE  
RECENT ADVANCES IN  
NATURAL LANGUAGE PROCESSING'2017

**PROCEEDINGS**

Varna, Bulgaria  
8 September 2017

ISBN 978-954-452-044-1

Designed and Printed by INCOMA Ltd.  
Shoumen, BULGARIA

## Preface

Biomedical NLP deals with the processing of healthcare-related text—clinical documents created by physicians and other healthcare providers at the point of care, scientific publications in the areas of biology and medicine, and consumer healthcare text such as social media blogs. Recent years have seen dramatic changes in the types and amount of data available to researchers in this field. Where most research on publications in the past has dealt with the abstracts of journal articles, we now have access to the full texts of journal articles via PubMedCentral. Where research on clinical documents has been hampered by a lack of availability of data, we now have access to large bodies of data through the auspices of the Cincinnati Children’s Hospital NLP Challenge, the i2b2 shared tasks ([www.i2b2.org](http://www.i2b2.org)), the TREC Electronic Medical Records track, Clinical TempEval series of tasks, the US-funded Strategic Health Advanced Research Projects Area 4 ([www.sharprn.org](http://www.sharprn.org)) and the Shared Annotated Resources (ShARE; <https://sites.google.com/site/shareclefehealth/taskdescription>; [www.clinicalnlpannotations.org](http://www.clinicalnlpannotations.org)) project. Meanwhile, the number of abstracts in PubMed continues to grow exponentially. Text in the form of blogs created by patients discussing various healthcare topics has emerged as another data source, with a new perspective on healthrelated issues. Connecting the information from the three main sources in multiple languages to the scientific community, the healthcare provider, and the healthcare consumer presents new challenges.

The Biomedical Natural Language Processing at RANLP 2017 provided a venue for presentations of current work in this field. The topics of papers presented at the workshop included information retrieval, part-of-speech tagging, multi-part knowledge frames population, extraction of numerical described values, resource creation, entity-centric information access, named entity recognition, confidence estimation for protein-protein relation discovery and association rule mining of clinical text and the biomedical literature.

The Workshop Organizers



## **The BioNLP Workshop associated with RANLP'17 is organised by:**

Svetla Boytcheva (IICT, Bulgarian Academy of Sciences)

Kevin Bretonnel Cohen (University of Colorado School of Medicine)

Guergana Savova (Harvard Medical School and Boston Children's Hospital)

Galia Angelova (IICT, Bulgarian Academy of Sciences)

## **The event is partially supported by:**

National Scientific Fund, Ministry of Education and Science, Bulgaria

## **Program Committee:**

Galia Angelova (Bulgarian Academy of Sciences)

Svetla Boytcheva (Bulgarian Academy of Sciences)

Kevin Cohen (U. Colorado School of Medicine)

Noa P. Cruz Diaz (Group of Research and Innovation in Biomedical Informatics, Biomedical Engineering and Health Economy. Institute of Biomedicine of Seville/ Virgen del Rocío University Hospital / CSIC / University of Seville)

George Giannakopoulos (NCSR Demokritos & SciFY NPC)

Agnieszka Mykowiecka (Polish Academy of Sciences)

Preslav Nakov (Qatar Computing Research Institute, HBKU)

Ivelina Nikolova (Bulgarian Academy of Sciences)

Georgios Petasis (NCSR "Demokritos")

Guergana Savova (Harvard Medical School and Boston Children's Hospital)

Frédérique Segond (Viseo Research)

Dimitar Tcharaktchiev (Medical University - Sofia)

## **Reviewers:**

Ekaterina L. Chernyak (Higher School of Economics, Moscow)

Dmitry Ilvovsky (Higher School of Economics, Moscow)

Natalia Korepanova (Higher School of Economics, Moscow)



## Table of Contents

<i>Document retrieval and question answering in medical documents. A large-scale corpus challenge.</i> Curea Eric .....	1
<i>Adapting the TTL Romanian POS Tagger to the Biomedical Domain</i> Maria Mitrofan and Radu Ion .....	8
<i>Discourse-Wide Extraction of Assay Frames from the Biological Literature</i> Dayne Freitag, Paul Kalmar and Eric Yeh .....	15
<i>Classification based extraction of numeric values from clinical narratives</i> Maximilian Zubke .....	24
<i>Understanding of unknown medical words</i> Natalia Grabar and Thierry Hamon .....	32
<i>Entity-Centric Information Access with Human in the Loop for the Biomedical Domain</i> Seid Muhie Yimam, Steffen Remus, Alexander Panchenko, Andreas Holzinger and Chris Biemann .....	42
<i>One model per entity: using hundreds of machine learning models to recognize and normalize biomedical names in text</i> Victor Bellon and Raul Rodriguez-Esteban .....	49
<i>Towards Confidence Estimation for Typed Protein-Protein Relation Extraction</i> Camilo Thorne and Roman Klinger .....	55
<i>Identification of Risk Factors in Clinical Texts through Association Rules</i> Svetla Boytcheva, Ivelina Nikolova, Galia Angelova and Zhivko Angelov .....	64
<i>POMELO: Medline corpus with manually annotated food-drug interactions</i> Thierry Hamon, Vincent Tabanou, Fleur Mouglin, Natalia Grabar and Frantz Thiessard .....	73
<i>Demo presentations:</i> <i>Annotation of Clinical Narratives in Bulgarian language</i> Ivaylo Radev, Kiril Simov, Galia Angelova and Svetla Boytcheva .....	81

