# SIXTH INTERNATIONAL CONFERENCE

## ATLB '24

# COMPUTATIONAL LINGUISTICS
# IN BULGARIA
# CLIB 2024

**9 – 10** September **2024**

**Sofia, Bulgaria**

# PROCEEDINGS

CLIB 2024 is organised by:

Department of Computational Linguistics
Institute for Bulgarian Language
Bulgarian Academy of Sciences

# PUBLICATION AND CATALOGUING INFORMATION

Proceedings of the

Sixth International Conference

COMPUTATIONAL LINGUISTICS IN BULGARIA



9 – 10 September 2024
Sofia, Bulgaria

## PROGRAMME COMMITTEE

of Linguistic Modelling and Knowledge Processing, Bulgarian Academy of Sciences, Bulgaria

**Kresimir Sojat** – University of Zagreb, Faculty of Humanities and Social Sciences, Croatia

**Ranka Stankovic** – University of Belgrade, Serbia

**Ivelina Stoyanova** – Institute for Bulgarian Language, Department of Computational Linguistics, Bulgarian Academy of Sciences, Bulgaria

**Stan Szpakowicz** – University of Ottawa, Canada

**Hristo Tanev** – Joint Research Centre of the European Commission, Italy

**Irina Temnikova** – Big Data for Smart Society Institute (GATE), Bulgaria

**Tinko Tinchev** – Sofia University, Faculty of Mathematics and Informatics, Bulgaria

**Maria Todorova** – Institute for Bulgarian Language, Department of Computational Linguistics, Bulgarian Academy of Sciences, Bulgaria

**Cristina Vertan** – University of Hamburg, Germany

**Shuly Wintner** – University of Haifa, Department of Computer Science, Israel

**Piek Vossen** – Free University of Amsterdam, the Netherlands

**Katerina Zdravkova** – University St Cyril and Methodius in Skopje, North Macedonia


## ORGANISING COMMITTEE

**Chair:**

**Svetlozara Leseva** – Institute for Bulgarian Language, Department of Computational Linguistics, Bulgarian Academy of Sciences, Bulgaria


**Atanas Atanasov** – Sofia University, Faculty of Slavic Studies, Bulgaria

**Rositsa Dekova** – Plovdiv University, Faculty of Philology, Department of English Studies, Bulgaria

**Nevena Grigorova** – Gate Institute, Sofia University, Bulgaria

**Dimitar Hristov** – Cleversoft, Bulgaria

**Nikolaos-Digenis Karagiannis** – Identrics, Bulgaria

**Hristina Kukova** – Institute for Bulgarian Language, Department of Computational Linguistics, Bulgarian Academy of Sciences, Bulgaria

**Todor Lazarov** – New Bulgarian University, Bulgaria

**Viktoria Petrova** – A1 Bulgaria, Bulgaria

**Valentina Stefanova** – Institute for Bulgarian Language, Department of Computational Linguistics, Bulgarian Academy of Sciences, Bulgaria

**Ekaterina Tarpomanova** – Sofia University, Faculty of Slavic Studies, Bulgaria

# Table of Contents

# PLENARY TALKS

## Large Language Models for the Real World: Explorations of Sparse, Cross-lingual Understanding and Instruction-Tuned LLMs

**Dr. Veselin Stoyanov (Tome AI, USA)**

---

Large language models (LLMs) have revolutionized NLP and the use of Natural Language in products. Nonetheless, there are challenges to the wide adoption of LLMs. In this talk, I will describe my explorations into addressing some of those challenges. I will cover work on sparse models addressing high computational costs, multilingual LLMs addressing the need to handle many languages, and work on instruction finetuning addressing the alignment between model outputs and human needs.

# Ten Years of Universal Dependencies

**Prof. Joakim Nivre (Uppsala University and RISE, Sweden)**

---

Universal Dependencies (UD) is a project developing cross-linguistically consistent treebank annotation for many languages, with the goal of facilitating multilingual parser development, cross-lingual learning, and parsing research from a language typology perspective. Since UD was launched almost ten years ago, it has grown into a large community effort involving over 500 researchers around the world, together producing treebanks for 148 languages and enabling new research directions in both NLP and linguistics. In this talk, I will review the history and development of UD and discuss challenges that we need to face when bringing UD into the future.

# Written Text Processing and the Adaptive Reading Hypothesis

**Prof. Vito Pirrelli (NRC, Institute for Computational Linguistics, Pisa, Italy)**

---

Oral reading requires the fine coordination of eye movements and articulatory movements. The eye provides access to the input stimuli needed for voice articulation to unfold at a relatively constant rate, while control on articulation provides internal feedback to oculomotor control for eye movements to be directed when and where a decoding problem arises.

A factor that makes coordination of the eye and the voice particularly hard to manage is their asynchrony. Eye movements are faster than voice articulation and are much freer to scan a written text forwards and backwards. As a result, given a certain time window, the eye can typically fixate more words than the voice can articulate.

According to most scholars, readers compensate for this functional asynchrony by using their phonological buffer, a working memory stack of limited temporal capacity where fixated words can be maintained temporarily, until they are read out loud. The capacity of the phonological buffer thus puts an upper limit on the distance between the position of the voice and the position of the eye during oral text reading, known as the eye-voice span.

In my talk, I will discuss recent reading evidence showing that the eye-voice span is the "elastic" outcome of an optimally adaptive viewing strategy, interactively modulated by individual reading skills and the lexical and structural features of a text. The voice span not only varies across readers depending on their rate of articulation, but it also varies within each reader, getting larger when a larger structural unit is processed. This suggests that skilled readers can optimally coordinate articulation and fixation times for text processing, adaptively using their phonological memory buffer to process linguistic structures of different size and complexity.