IWCLUL 2025


**The 10th International Workshop on Computational Linguistics for Uralic Languages**


**Proceedings of the Workshop**


December 10-12, 2025

Order copies of this and other ACL proceedings from:

# Preface

Welcome to the Proceedings of the 10th International Workshop on Computational Linguistics for Uralic Languages (IWCLUL), a dedicated venue for scholars, practitioners, and researchers working in computational linguistics with a particular emphasis on Uralic languages. This year's workshop continues the IWCLUL tradition of fostering interdisciplinary exchange, shared learning, and a collective dedication to tackling the linguistic, technological, and social issues surrounding Uralic languages in the digital era. The Uralic language family—stretching across Europe and Asia and including languages as varied as Finnish, Hungarian, and the endangered Udmurt and Khanty—brings with it a distinct set of computational challenges. These languages often feature rich morphology, agglutinative patterns, and unique syntactic and phonological structures, all of which demand specialized methods in computational analysis and language technology. Our workshop aims to highlight these complexities and encourage the creation of innovative tools and approaches that not only facilitate digital use of these languages but also contribute to their long-term vitality.

This year, IWCLUL saw a strong and wide-ranging set of submissions from the international research community, demonstrating the sustained interest in computational work on Uralic languages. The accepted papers span numerous topics and language varieties, offering insights into both well-documented and lesser-resourced Uralic languages. This diversity reflects the ongoing growth of the field and the many directions in which researchers are advancing computational linguistics for Uralic languages.

We are also pleased to celebrate Jack Rueter, who was honored with a lifetime achievement award by ACL SIGUR for his decades-long dedication to Uralic language research, documentation, and technology development. His work has had a lasting impact on the community and serves as an inspiration for future generations of researchers.

We hope that these proceedings inspire continued research and collaboration in computational linguistics for Uralic languages. May the insights, methodologies, and resources shared here contribute to meaningful advances in the field and foster an inclusive future for Uralic languages in the digital landscape.

Sincerely, The IWCLUL 2025 Organizing Committee

# Organizing Committee

**Organizers**

Mika Hämäläinen, Metropolia University of Applied Sciences
Flammie Pirinen, Arctic University of Norway
Lev Kharlashkin, Metropolia University of Applied Sciences
Eiaki V. Morooka, Metropolia University of Applied Sciences
Michael Rießler, University of Eastern Finland
Maarit Koponen, University of Eastern Finland
Ilia Moshnikov, University of Eastern Finland
Diana Kulashekhar, University of Eastern Finland
Mahla Baniasadi, University of Eastern Finland
Nurstrat prova, University of Eastern Finland
Varoon Bakshi, University of Eastern Finland

# Program Committee

Laszlo Fejes, Hungarian Research Centre for Linguistics

Gunta Klava, University of Latvia
Sourav Das, Indian Institute of Information Technology Kalyani
Balazs Indig, Eotvos Lorand University
Laleh Davoodi, Abo Akademi University
Flammie A Pirinen, Norgga arktalas universitehta
Trond Trosterud, University of Tromso
Csilla Horvath, University of Szeged
Iaroslav Chelombitko, Neapolis University Paphos
Aleksei Dorkin, University of Tartu
Sebastian Oliver Eck, University of Oxford
Khalid Alnajjar, F-Secure Oyj
Jules Bouton, Universite Paris Cite
Valts Ernstreits, University of Latvia
Michael Riessler, University of Eastern Finland
Timofey Arkhangelskiy, Universitat Hamburg
Samy Ouzerrout, Universite d'Orleans
Jyoti Kunal Shah, Automatic Data Processing
Youngsook Song, Lablup
David Dale, FAIR at Meta
Abhi Desai, New England College

# Table of Contents

# Program

**Wednesday, December 10, 2025**

14:00 - 18:00      *Karelian Workshop*

**Thursday, December 11, 2025**

**Friday, December 12, 2025**

10:00 - 10:00    *Day 3 Opening*

10:00 - 11:00    *Oral Session 4*

*Benchmarking Large Language Models for Lemmatization and Translation of Finnic Runosongs*
Lidia Pivovarova, Kati Kallio, Antti Kanner, Jakob Lindström, Eetu Mäkelä, Liina Saarlo, Kaarel Veskis and Mari Väina

*Fine-Tuning Whisper for Kildin Sami*
Enzo Gamboni

*Digitization Work at the Finno-Ugrian Society: Livonian Case Study*
Niko Partanen, Jack Rueter and Valts Ernštreits

11:00 - 12:00    *Lunch*

12:00 - 13:00    *Oral Session 5*

*Siberian Ingrian Finnish: FST and IGTs*
Ivan Ubaleht

*Case–Number Dissociation in Finnish Noun Embeddings:fastText vs. BERT Layer Effects*
Alexandre Nikolaev, Yu-Ying Chuang and R. Harald Baayen

*Evaluating OpenAI GPT Models for Translation of Endangered UralicLanguages: A Comparison of Reasoning and Non-Reasoning Architectures*
Yehor Tereschenko, Mika Hämäläinen and Svitlana Myroniuk

13:00 - 13:20    *Coffee Break*

13:20 - 14:20    *SIGUR Business Meeting*