# NEALT

# Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)

May 22 - 24, 2023

Tórshavn, Faroe Islands

Editors: Tanel Alumäe and Mark Fishel

NoDaLiDa 2023

**24th Nordic Conference on Computational Linguistics (NoDaLiDa)**

**Proceedings of the Conference**

May 22-24, 2023

The NoDaLiDa organizers gratefully acknowledge the support from the following sponsors.

**Silver**



**Bronze**



**Other**

Volume Editors:
Tanel Alumäe and Mark Fishel

# Message from the General Chair

It is my great pleasure and honor to welcome you to the 24th Nordic Conference on Computational Linguistics (NoDaLiDa 2023)!

After a couple of years' worth of conferences cancelled or held online (including the previous NoDaLiDa) we are extremely happy that NoDaLiDa 2023 is an onsite event. This is especially exciting given that for the first time in the history of NoDaLiDa conferences it takes place in Tórshavn, Faroe Islands.

The conference features three types of papers: long, short and demo papers. We are truly grateful to all the authors of papers submitted to this year's conference, with 130 papers submitted, a more than 40% increase over last year's yield! In total, we accepted 79 papers: 49 long papers, 26 short papers and 4 demo papers. More than half of the accepted papers are student papers, in which the first author is a student (29 long, 17 short and 2 demo papers). We would like to thank the 113 members of the program committee who reviewed the papers for their contributions!

The 79 accepted papers are grouped into 12 oral and 2 poster sessions. In addition to these regular sessions the conference program also includes three keynote talks. We would like to extend our gratitude to the keynote speakers for agreeing to present their work at NoDaLiDa. Georg Rehm from DFKI will talk on the topic of "Towards Digital Language Equality in Europe: An Overview of Recent Developments". Hjalmar P. Petersen will talk about "Aspects of the structure of Faroese". Marta R. Costa-Jussà from Meta will talk about "No-language-left-behind: Scaling Human-Centered Machine Translation and Toxicity at Scale".

The main conference program is preceded by three workshops: NLP for Computer-Assisted Language Learning (NLP4CALL), the Constraint Grammar Workshop and Resources and representations for under-resourced languages and domains (RESOURCEFUL'2023). We thank the workshop organizers for their efforts and for expanding the main conference program with a focus on more specific research topics.

I would like to express sincere gratitude to the entire team behind organizing NoDaLiDa 2023. I was honored to receive the invitation to serve as the general chair from the NEALT board; thank you for trusting me with this role. My deepest gratitude goes to Tanel Alumäe for serving as the publications chair and his active participation, Inguna Skadiņa for serving as the workshop chair as well as Iben Nyholm Debess for serving as the main local chair and smoothly handling all associated aspects of conference organization. I also want to thank the rest of the program chairs, Lilja Øvrelid and Christian Hardmeier and the local co-chairs Bergur Djurhuus Hansen, Peter Juel Henrichsen and Sandra Saxov Lamhauge. Thank you everyone for your contributions, you are awesome!

NoDaLiDa 2023 received financial support from several institutions and we would like to thank them here: NEALT, Dictus, Málráðið, Tórshavnar kommuna, BankNordik, Digitaliseringsstyrelsen, University of the Faroe Islands, Nationella språkbanken, Elektron and Formula.

Welcome and enjoy the 24th Nordic Conference on Computational Linguistics!

Mark Fishel, General Chair

Tartu

May 2023

# Organizing Committee

**General Chair**

    Mark Fishel, University of Tartu

**Program Chairs**

    Tanel Alumäe, Tallinn University of Technology (Publication Chair)
    Inguna Skadina, University of Latvia (Workshop Chair)
    Christian Hardmeier, IT University of Copenhagen
    Lilja Øvrelid, University of Oslo

**Local Chair**

    Iben Nyholm Debess, University of the Faroe Islands

**Local Co-Chairs**

    Bergur Djurhuus Hansen, University of the Faroe Islands
    Peter Juel Henrichsen, Danish Language Council
    Sandra Saxov Lamhauge, University of the Faroe Islands

## Reviewers

# Invited Talk: Towards Digital Language Equality in Europe: An Overview of Recent Developments

**Georg Rehm**
German Research Center for Artifical Intelligence

Digital Language Equality (DLE) "is the state of affairs in which all languages have the technological support and situational context necessary for them to continue to exist and to prosper as living languages in the digital age", as we specified in one of our key reports of the EU project European Language Equality (ELE). Our empirical findings suggest that Europe is currently very far from having a situation in which all our languages are supported equally well through technologies. In this presentation, I'll give an overview of the two ELE projects and their main results and findings with a special focus on the Nordic languages (including insights from the FSTP projects supported through ELE2). This will also include a brief look back into the past, especially discussing the question if and where we have seen progress in the last, say, 15 years. Furthermore, I'll present an overview of our main strategic recommendations towards the European Union in terms of bringing about DLE in Europe by 2030. The presentation will conclude with a look at other relevant activities in Europe, including, critically the Common European Language Data Space project, which started in early 2023.

# Invited Talk: No-language-left-behind: Scaling Human-Centered Machine Translation and Toxicity at Scale

**Marta R. Costa-jussà**
Meta AI

Machine Translation systems can produce different types of errors, some of which are characterized as critical or catastrophic due to the specific negative impact that they can have on users. In this talk, we focus on one type of critical error: added toxicity. We evaluate and analyze added toxicity in the context of NLLB-200 that open-sources models capable of delivering evaluated, high-quality translations directly between 200 languages. An automatic toxicity evaluation shows that added toxicity across languages varies from 0% to 5%. The output languages with the most added toxicity tend to be low-resource ones, and the demographic axes with the most added toxicity include sexual orientation, gender and sex, and ability. Making use of the input attributions allows us to further explain toxicity and our recommendations to reduce added toxicity are to curate training data to avoid mistranslations, mitigate hallucination and check unstable translations.

# Invited Talk: Aspects of the Structure of Faroese

**Hjalmar P. Petersen**
University of Faroe Islands

Phonological changes and later morphologization have led to different complex alternations in Faroese. These are argued to emerge especially in small languages, with little contact and tight networks. The alternations will be exemplified with 'skerping', palatalization, glide insertion and the quantity-shift. There will be a discussion of the morphology-phonology interface, where the suggestion is that Faroese has 3 strata, stem1, stem2 and a word- strata. Syntactic variation and different construction will be addressed and illustrated; in this context reflexives are included and the present reorganization of the case system of complements of prepositions, where speakers use semantic and structural case in a certain way.

# Table of Contents

xi