

VarDial 2025

**VarDial 2025 - The Twelfth Workshop on NLP for Similar
Languages, Varieties and Dialects**

Proceedings of the Workshop

January 19, 2025

©2025 The Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)
317 Sidney Baker St. S
Suite 400 - 134
Kerrville, TX 78028
USA
Tel: +1-855-225-1962
acl@aclweb.org

ISBN 979-8-89176-208-4

Preface

These proceedings include the 17 papers presented at the Twelfth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial 2025), co-located with the 31st International Conference on Computational Linguistics (COLING 2025). VarDial was held in Abu Dhabi, UAE.

Despite the short interval between the 2024 and 2025 editions of VarDial, we are glad to see that VarDial continues to serve the community as the main venue for researchers interested in the computational processing of language variation. The papers accepted this year address a wide range of topics, such as normalization and dialectal translation, native language identification, and slot and intent detection. We also see several papers making use of and evaluating large language models on variety-related tasks. Once again, these proceedings are characterized by great linguistic diversity, with work on regional English dialects, Portuguese, Luxemburgish, Church Slavic, and Arabic, to name just a few.

As in previous editions, VarDial 2025 features an evaluation campaign with the NorSID shared task on slot, intent and dialect identification for Norwegian dialects. Slot and intent detection were already included in a VarDial shared task in 2023, but without including Norwegian data. Likewise, language and dialect identification tasks have been very common at past editions of VarDial, but this is the first dialect identification featuring varieties of Norwegian. This volume includes the system description papers prepared by the four participating teams, as well as a report written by the task organizers summarizing the results and findings of the evaluation campaign.

Finally, we would like to take this opportunity to thank all the shared task organizers and the participants for their hard work. We further thank the VarDial program committee members for being an important part of the workshop's success.

The VarDial workshop organizers:

Yves Scherrer, Tommi Jauhiainen, Nikola Ljubešić, Preslav Nakov, Jörg Tiedemann, and Marcos Zampieri

<http://sites.google.com/view/vardial-2025/>

Organizing Committee

Organizers:

Yves Scherrer, University of Oslo (Norway)

Tommi Jauhiainen, University of Helsinki (Finland)

Nikola Ljubešić, Jožef Stefan Institute and University of Ljubljana (Slovenia)

Preslav Nakov, Mohamed bin Zayed University of Artificial Intelligence (UAE)

Jörg Tiedemann, University of Helsinki (Finland)

Marcos Zampieri, George Mason University (USA)

Program Committee

Program Committee:

Noëmi Aepli (University of Zurich, Switzerland)
César Aguilar (Universidad Veracruzana, Mexico)
Sina Ahmadi (George Mason University, United States)
Laura Alonso y Alemany (University of Cordoba, Argentina)
Jorge Baptista (University of Algarve and INESC-ID, Portugal)
Delphine Bernhard (University of Strasbourg, France)
Gabriel Bernier-Colborne (National Research Council, Canada)
Verena Blaschke (LMU Munich, Germany)
Francis Bond (Nanyang Technological University, Singapore)
David Chiang (University of Notre Dame, United States)
Steven Coats (University of Oulu, Finland)
Çağrı Çöltekin (University of Tübingen, Germany)
Stefanie Dipper (Ruhr University Bochum, Germany)
Sascha Diwersy (University of Montpellier, France)
Mark Dras (Macquarie University, Australia)
Jonathan Dunn (University of Illinois Urbana-Champaign, United States)
Pablo Gamallo (University of Santiago de Compostela, Spain)
Rob van der Goot (IT University Copenhagen, Denmark)
Cyril Goutte (National Research Council, Canada)
Nizar Habash (New York University Abu Dhabi, UAE)
Radu Ionescu (University of Bucharest, Romania)
Anjali Kantharuban (Carnegie Mellon University, United States)
John McCrae (University of Galway, Ireland)
Surafel Melaku Lakew (FBK , Italy)
Aleksandra Miletic (University of Helsinki, Finland)
Filip Miletic (University of Stuttgart, Germany)
John Nerbonne (University of Groningen, Netherlands and University of Freiburg, Germany)
Ekaterina Lapshinova-Koltunski (University of Hildesheim, Germany)
Lung-Hao Lee (National Yang Ming Chiao Tung University, Taiwan)
Maciej Ogrodniczuk (Institute of Computer Science, Polish Academy of Sciences, Poland)
Petya Osenova (Bulgarian Academy of Sciences, Bulgaria)
Jelena Prokic (Leiden University, Netherlands)
Christoph Purschke (University of Luxembourg, Luxembourg)
Francisco Rangel (Autoritas Consulting, Spain)
Reinhard Rapp (University of Mainz, Germany)
Tanja Samardžić (University of Zurich, Switzerland)
Serge Sharoff (University of Leeds, United Kingdom)
Miikka Silfverberg (University of British Columbia, Canada)
Milena Slavcheva (Bulgarian Academy of Sciences, Bulgaria)
Aarohi Srivastava (University of Notre Dame, United States)
Marco Tadić (University of Zagreb, Croatia)
Joel Tetreault (Dataminr, United States)
Pidong Wang (Google Inc., United States)
Taro Watanabe (Google Inc., Japan)

Table of Contents

<i>Findings of the VarDial Evaluation Campaign 2025: The NorSID Shared Task on Norwegian Slot, Intent and Dialect Identification</i>	
Yves Scherrer, Rob van der Goot and Petter Mæhlum	1
<i>Information Theory and Linguistic Variation: A Study of Brazilian and European Portuguese</i>	
Diego Alves	9
<i>Leveraging Open-Source Large Language Models for Native Language Identification</i>	
Yee Man Ng and Ilia Markov	20
<i>Adapting Whisper for Regional Dialects: Enhancing Public Services for Vulnerable Populations in the United Kingdom</i>	
Melissa Torgbi, Andrew Clayman, Jordan J. Speight and Harish Tayyar Madabushi	29
<i>Large Language Models as a Normalizer for Transliteration and Dialectal Translation</i>	
Md Mahfuz Ibn Alam and Antonios Anastasopoulos	39
<i>Testing the Boundaries of LLMs: Dialectal and Language-Variety Tasks</i>	
Fahim Faisal and Antonios Anastasopoulos	68
<i>Text Generation Models for Luxembourgish with Limited Data: A Balanced Multilingual Strategy</i>	
Alistair Plum, Tharindu Ranasinghe and Christoph Purschke	93
<i>Retrieval of Parallelizable Texts Across Church Slavic Variants</i>	
Piroska Lendvai, Uwe Reichel, Anna Jouravel, Achim Rabus and Elena Renje	105
<i>Neural Text Normalization for Luxembourgish Using Real-Life Variation Data</i>	
Anne-Marie Lutgen, Alistair Plum, Christoph Purschke and Barbara Plank	115
<i>Improving Dialectal Slot and Intent Detection with Auxiliary Tasks: A Multi-Dialectal Bavarian Case Study</i>	
Xaver Maria Krückl, Verena Blaschke and Barbara Plank	128
<i>Regional Distribution of the /eI/-/æI/ Merger in Australian English</i>	
Steven Coats, Chloé Diskin-Holdaway and Debbie Loakes	147
<i>Learning Cross-Dialectal Morphophonology with Syllable Structure Constraints</i>	
Salam Khalifa, Abdelrahim Qaddoumi, Jordan Kodner and Owen Rambow	157
<i>Common Ground, Diverse Roots: The Difficulty of Classifying Common Examples in Spanish Varieties</i>	
Javier A. Lopetegui, Arij Riabi and Djamé Seddah	168
<i>Add Noise, Tasks, or Layers? MaiNLP at the VarDial 2025 Shared Task on Norwegian Dialectal Slot and Intent Detection</i>	
Verena Blaschke, Felicia Körner and Barbara Plank	182
<i>LTG at VarDial 2025 NorSID: More and Better Training Data for Slot and Intent Detection</i>	
Marthe Midtgaard, Petter Mæhlum and Yves Scherrer	200

HiTZ at VarDial 2025 NorSID: Overcoming Data Scarcity with Language Transfer and Automatic Data Annotation

Jaione Bengoetxea, Mikel Zubillaga, Ekhi Azurmendi, Maite Heredia, Julen Etxaniz, Markel Ferro and Jeremy Barnes 209

CUFE@VarDial 2025 NorSID: Multilingual BERT for Norwegian Dialect Identification and Intent Detection

Michael Ibrahim 220