NAACL 2025

**Annual Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics**

**Proceedings of the Conference
Volume 1: Long Papers**

April 29 - May 4, 2025

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

**Platinum**





**Gold**



**Bronze**

# Diversity and Inclusion Ally

# Message from the General Chair

**Message from the General Chair**  Welcome to the 2025 meeting of the Nations of the Americas Chapter of the Association for Computational Linguistics! I am proud to help organize the first NAACL conference to carry the new name of our organization, one that emphasizes inclusion for all of the Americas. I am also pleased to welcome you to Albuquerque, New Mexico, a state whose unique blend of cultural influences will make for an excellent backdrop for NAACL 2025, especially with this year's special theme on NLP in a Multicultural World.

This year's program benefited from a now mature ACL rolling review process. I would like to extend a big thank you to our ARR Editors in Chief, Viviane Moreira, Anna Rogers, and Michael White, who were very helpful not just with reviewing for the main conference, but who also shared their OpenReview expertise with chairs from our other tracks. We also benefited from the helpful advice of last year's NAACL general and program chairs: Katrin Erk, Kevin Duh, Helena Gomez, and Steven Bethard. Finally, Ryan Cotterell stepped up to help with the publications process and software, even while not serving as publications chair.

Of course, a conference of this magnitude cannot come together without some drama; in our case, we had some unexpected funding shortages from traditional government sources. We would like to extend a huge thank you to the boards of both the ACL and NAACL for filling those funding gaps, ensuring that our important D&I, volunteer, and student author support programs continue to thrive.

The virtual component of our conference is crucial to an inclusive, affordable experience for all NAACL authors and attendees. This year, we made small refinements to the hybrid format, mirroring NAACL 2024's choices to combine a virtual poster session via Gather with asynchronous online content via Underline. Our virtual poster session will be hosted on May 6, the Tuesday after the conference, in the hopes that promoting it at the conference's plenary sessions will help boost attendance. We opted not to have virtual oral presentations at the in-person conference, as those continue to be tricky to get right. To participants, virtual as well as in-person: Please let us know what worked for you and what did not, so we can continue to improve the hybrid experience.

The job of General Chair is a strange one, as it mostly involves cheering on many other people as they do amazing work. I have been fortunate to have been teamed with an excellent set of program chairs; to Luis Chiruzzo, Alan Ritter, and Lu Wang: thanks for everything, I'm very proud of what we've built together. I'd also like to extend my heartfelt thanks to Jenn Rachford (ACL) and Damira Mršic (Underline) who provide the knowledge, continuity and professionalism to bring all of this together.

Many thanks also to:

- Workshop chairs: Saab Mansour, Kenton Murray, and Alexis Palmer

- Tutorial chairs: Maria Lomeli, Swabha Swayamdipta, and Rui Zhang

- Demo chairs: Nouha Dziri, Shizhe Diao, and Sean (Xiang) Ren

- Industry track chairs: Weizhu Chen, Xue-Yong Fu, Mohammad Kachuee, and Yi Yang

- Student research workshop chairs: Abteen Ebrahimi, Emmy Liu, and Samar Haider, and their faculty advisors Maria Leonor Pacheco and Shira Wein

- Publication chairs: Arman Cohan, Manling Li, and Yichao Zhou

- Website chairs: Arya McCarthy and Vered Shwartz

- Publicity and social media chairs: Eleftheria Briakou, Tuhin Chakrabarty, and Ximena Gutierrez-Vasques

- Diversity and inclusion chairs: Akiko I. Eriguchi, Chi-Kiu (Jackie) Lo, and Niloofar Mireshghallah

- Sponsorship chairs: Prithviraj (Raj) Ammanabrolu and Maha Elbayad

- Volunteers chairs: Robin Jia and David Mortensen

- Ethics chairs: Manuel Mager and Yulia Tsvetkov

- Handbook chair: Winston Wu

- Best paper committee chairs: Marine Carpuat and Anna Rumshisky

- Visa chairs: Eduardo Blanco and Parisa Kordjamshidi

- Virtual infrastructure chair: Jieyu Zhao

Whenever possible, I tried to populate each committee with someone who had served in the same role in NAACL 2024, to provide continuity, so I'll extend an extra thanks to all chairs who accepted this second year of service. Thanks also to the members of the ACL and NAACL Executive Committees for their support, feedback, and advice.

Finally, I would like to thank all authors, invited speakers and panelists, area chairs and reviewers, volunteers and session chairs, and all attendees, in-person and virtual. The conference is nothing without you.

Welcome again and enjoy the conference!

**Colin Cherry**
Google
NAACL 2025 General Chair

# Message from the Program Chairs

**Message from the Program Chairs**  Welcome to the 2025 Annual Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics! NAACL 2025 marks an exciting milestone as the first conference held under our new name, reflecting our commitment to greater inclusivity across the diverse communities of the Americas. NAACL 2025 is a hybrid conference, and we are excited to have attendees and presenters join us both in person in Albuquerque and online from around the globe. We are thrilled to welcome you to what promises to be a vibrant and engaging conference.

**Special Theme: NLP in a Multicultural World**  Current NLP tools and models, particularly large language models (LLMs), rely on vast amounts of data for training. However, this data often over-represents a small number of dominant languages, and even within those, tends to prioritize certain geographical or cultural varieties. As a result, a long tail of under-represented languages, dialects, and cultural contexts remains largely overlooked by the NLP community. For NAACL 2025, we introduced a special theme track on "NLP in a Multicultural World." With this theme track, we sought to foster discussion and research on how NLP can better serve the linguistic and cultural diversity of the world. We encouraged contributions on topics such as cultural localization of language models, new NLP applications to support people from diverse cultures, revitalization or refunctionalization of endangered or sleeping languages, analysis of cultural biases in language models, and historical considerations and diachronic analysis. This track was dedicated to developing more inclusive, culturally aware NLP techniques that reflect and support the vibrant multicultural world we live in.

We received 71 submissions to the special theme, of which 23 were accepted for presentation at the conference. We hope these papers spark meaningful conversations and inspire future work in this important and evolving area of research.

**Review Process**  3,185 papers were submitted to the October ARR cycle, of which we estimate 3,099 were intended to be submitted to NAACL based on the "preferred venue" field in the submission form. We also received 147 papers from previous ARR cycles committed to NAACL. The program chairs recruited 98 Senior Area Chairs to view reviews and metareviews provided by ARR and make final recommendations on which papers to accept to both the main conference and Findings. 1,432 Area Chairs wrote metareviews for ARR, and 10,648 reviewers wrote reviews for the submitted papers.

**Acceptance Rate**  Calculating an acceptance rate is challenging due to the multi-step ARR review process, in which papers are first submitted to ARR to get reviews, then authors commit their papers (together with reviews) to a specific *ACL conference. Of the 3,185 papers submitted to the October ARR cycle, we estimate that 3,099 intended to submit to NAACL. Based on this information, we estimate that 22% of papers submitted to the October cycle and intended for NAACL were accepted to the main conference, and another 15% were accepted to Findings, bringing the total estimated acceptance rate for papers accepted to be presented at the conference (Main + Findings) to 37%. Out of the 1,647 papers committed to NAACL with ARR reviews, 719 were accepted to the main conference, and 477 were accepted to Findings. 40 papers were desk rejected or withdrawn.

**Presentation Format**  At NAACL 2025, papers were assigned one of three possible presentation modes: in-person participants could be assigned oral or poster presentations, while virtual participants could present posters. We selected 246 of the papers accepted to the main conference as oral presentations, and

the rest of them were assigned as posters, together with all the Findings papers. Oral presentations were assigned a 15-minute slot, with 12 minutes for presentation and 3 minutes for questions. For choosing papers as oral presentations, we first split all papers according to track, sorted them according to overall score (considering review, metareview, SAC recommendation), and took into consideration the authors' presentation preference. Then we grouped papers in sets of 6. Some tracks had very few accepted papers, so some of them were grouped together to form areas of affinity.

**Program Format**   NAACL 2025 consists both of an in-person and a virtual conference, held on different days. The virtual part of the conference is held after the in-person one and a few days later (on May 6), so participants traveling home after the in-person conference could attend the virtual conference. The conference program includes 3 keynote speakers: Rada Mihalcea (University of Michigan), Mike Lewis (Meta), and Josh Tenenbaum (Massachusetts Institute of Technology). 260 papers are scheduled to be presented as oral presentations (also including papers from TACL, CL, and the industry track), 594 papers are scheduled as in-person posters, and 256 virtual posters.

**Gratitude**   NAACL 2025 would not have been possible without the hard work of all people involved. We thank everyone who contributed, including:

- The General Chair, Colin Cherry.

- The ARR Editors-in-Chief of the October 2024 cycle: Viviane Moreira, Anna Rogers, Michael White.

- The OpenReview team, especially Rachel Smart.

- The 98 Senior Area Chairs.

- The 1,432 Area Chairs and 10,648 Reviewers.

- The best paper committee chairs, Marine Carpuat and Anna Rumshisky.

- The ethics chairs, Yulia Tsvetkov and Manuel Mager, and their team of reviewers.

- The website chairs, Vered Shwartz and Arya McCarthy.

- The publication chairs, Yichao Zhou, Manling Li, and Arman Cohan.

- The publicity chairs, Ximena Gutierrez-Vasques, Eleftheria Briakou, and Tuhin Chakrabarty.

- The volunteers chairs, Robin Jia and David Mortensen.

- The visa chairs, Eduardo Blanco and Parisa Kordjamshidi.

- The ACL Anthology Director, Matt Post, and his team.

- The Program Chairs of EMNLP 2024 (Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung (Vivian) Chen) and NAACL 2024 (Kevin Duh, Helena Gomez, and Steven Bethard).

- Damira Mršic and the Underline Team.

- Jenn Rachhford and the entire conference support staff.

**Luis Chiruzzo**, Universidad de la República, Uruguay
**Alan Ritter**, Georgia Institute of Technology
**Lu Wang**, University of Michigan
NAACL 2025 Program Committee Co-Chairs
April 2025

# Table of Contents

xiii

xxii

xxxi