MRL 2024

# The 4th Workshop on Multilingual Representation Learning

# Proceedings of the Workshop

November 16, 2024

# Organizing Committee

**Organizers**

David Ifeoluwa Adelani, McGill University, Canada
Duygu Ataman, New York University, USA
Mammad Hajili, Microsoft, USA
Raghav Mantri, New York University, USA
David Stap, University of Amsterdam, Netherlands
Jonne Sälevä, Brandeis University, USA
Francesco Tinner, University of Amsterdam, Netherlands
Abraham Owodunni, Ohio State University, USA

# Program Committee

**Reviewers**

David Ifeoluwa Adelani, Manuel Bobie Amankwatia, Catherine Arnett

Travis M. Bartley, Vishal Bhalla

Jiajing Chen, Xupeng Chen, Olutosoye Taiwo Christian, Zoltan Csaki

Konstantin Dobler, Koel Dutta Chowdhury

Yassine El Kheir, Abdellah El Mekki

Senkang Hu

Yusif Ibrahimov, Mironshoh Inomjonov, Jafar Isbarov

Ainaz Jamshidi, Jiby Mariya Jose

Yixiao Kang, Zhengjian Kang, Christopher Klamm, Hongzhi Kuai

Senyu Li, Yueqian Lin, Weisi Liu

Pranita Yogesh Mahajan, Nimshi Venkat Meripo

Muhammad Amin Nadim, Usman Nawaz

Esther Odunayo Oduntan, Peter Oseghale Ohue, Yewande Ojo

Chester Palen-Michel

Shaibal Saha, Shubham Shukla, Janet Yunchen Sung, Jonne Sälevä

Shaomu Tan, Wenjia Tan, Yihang Tao, Shailja Thakur

Vajratiya Vajrobol, Deepali Verma

Sahil Walia, Azmine Toushik Wasi, Di Wu

Sadia Zaib, Zhehao Zhang, Xufeng Zhao, Huichi Zhou, Ziqi Zhou

# Keynote Talk
# Invited talk 1

**Karen Livescu**
TTI at Chicago
**2024-11-16 9:10** –

**Bio:** Karen Livescu is a Professor at the Toyota Technological Institute at Chicago (TTIC). Her research focuses on speech and language processing and related areas of machine learning. She obtained her PhD in Computer Science MIT 2005, working in the Spoken Language Systems group of the Computer Science and Artificial Intelligence Laboratory.

# Keynote Talk
# Invited talk 2

**Hila Gonen**
University of Washington
**2024-11-16 9:50** –

**Bio:** Hila Gonen is a postdoctoral Researcher at the Paul G. Allen School of Computer Science & Engineering at the University of Washington working on Natural Language Processing. Her research focuses on two goals: (1) making cutting-edge language technology available and fair across speakers of different languages and users of different socio-demographic groups; (2) developing algorithms and methods for controlling the model's behavior. Prior to joining UW, she completed a Ph.D in Computer Science at the NLP lab at Bar Ilan University.

# Keynote Talk
# Invited talk 3

**Sebastian Ruder**
Cohere for AI
**2024-11-16 16:00** –

**Bio:** Sebastian Ruder is a research scientist at Cohere based in Berlin, Germany working on making large language models (LLMs) multilingual. He completed his PhD in Natural Language Processing and Deep Learning at the Insight Research Centre for Data Analytics.

# Table of Contents

# Program

**Saturday, November 16, 2024**

09:00 - 09:10  *Opening Remarks*

09:10 - 09:50  *Invited Talk by Karen Livescu*

09:50 - 10:30  *Invited Talk by Hila Gonen*

10:30 - 11:00  *Coffee Break*

11:00 - 12:30  *Poster Session*

12:30 - 14:00  *Lunch Break*

14:00 - 14:15  *Findings Paper*

14:15 - 14:30  *Winning Team Presentation*

14:30 - 15:00  *Best Paper*

15:00 - 15:30  *Honorable Mentions*

15:30 - 16:00  *Coffee Break*

16:00 - 16:50  *Invited Talk by Sebastian Ruder*

16:50 - 17:00  *Closing Remarks*