ALVR 2024

**The 3rd Workshop on Advances in Language and Vision Research**

**Proceedings of the Workshop**

August 16, 2024

# Introduction

Welcome to the 3rd Workshop on Advances in Language and Vision Research. Co-located with ACL 2024, the workshop is scheduled for August 16, 2024. To facilitate the participation of the global NLP and CV community, we continue running the workshop in a hybrid format.

Language and vision research has attracted great attention from both natural language processing (NLP) and computer vision (CV) researchers. Gradually, this area is shifting from passive perception, templated language, and synthetic imagery/environments to active perception, natural language, and photo-realistic simulation or real-world deployment. The workshop covers (but is not limited to) the following topics:

- Self-supervised vision and language pre-training;

- New tasks and datasets that provide real-world solutions in language and vision;

- Text-to-image/video generation and text-guided image/video editing;

- External knowledge integration in visual and language understanding;

- Visually-grounded natural language understanding and generation;

- Language-grounded visual recognition and reasoning;

- Language-grounded embodied agents, e.g., vision-and-language navigation;

- Visually-grounded multilingual study, e.g., multimodal machine translation;

- Shortcomings of the existing large vision & language models on downstream tasks and solutions;

- Ethics and bias in large vision & language models;

- Multidisciplinary study that may involve linguistics, cognitive science, robotics, etc.;

- Explainability and interpretability in large vision & language models.

Our agenda features keynote speeches, hybrid talk sessions both for long and short papers, and poster sessions. This year we received 35 submissions, and after a thorough peer-review process, 31 papers were accepted. Among the accepted papers, 18 are archive papers and 13 are non-archive papers.

We would like to deeply thank all the authors, committee members, keynote speakers, and participants for helping us grow this research community both in quantity and quality.

Workshop Chairs

Jing Gu, UC Santa Cruz
Tsu-Jui Fu, UC Santa Barbara
Drew Hudson, Google DeepMind
Asli Celikyilmaz, Fundamentals AI Research @ Meta
William Wang, UC Santa Barbara

# Organizing Committee

**General Chair**

Jing Gu, University of California Santa Cruz, USA
Tsu-Jui (Ray) Fu, Apple, USA
Drew Hudson, Google DeepMind, USA
Asli Celikyilmaz, Fundamentals AI Research (FAIR) @ Meta, USA
William Wang, University of California Santa Barbara, USA

# Program Committee

Peiyan Zhang, Hong Kong University of Science and Technology
Siqiao Zhao, Morgan Stanley
Kaizhi Zheng, University of California, Santa Cruz
Chang Zhou, Columbia University
Wanrong Zhu, University of California, Santa Barbara
Fangrui Zhu, Northeastern University

# Table of Contents

# Program

**Friday, August 16, 2024**

09:00 - 09:15    *Welcome Speech*

09:15 - 10:15    *Invited Talk by Dr. Alane Suhr*

10:15 - 10:30    *Poster Session 1*

10:30 - 11:00    *Coffee Break*

11:00 - 12:00    *Invited Talk by Dr. Angel Chang*

12:00 - 13:00    *Invited Talk by Dr. Daniel Fried*

13:00 - 13:45    *Lunch Break*

13:45 - 14:45    *Invited Talk by Roozbeh Mottaghi*

14:45 - 15:45    *Invited Talk by Dr. Heng Ji*

15:45 - 16:15    *Coffee Break*

16:15 - 17:15    *Invited Talk by Xin (Eric) Wang*

17:15 - 17:45    *Poster Session 2*

17:45 - 18:00    *Closing Remarks*