

EvalMG 2025

The Workshop of Evaluation of Multi-Modal Generation

**Proceedings of the The First Workshop of Evaluation of
Multi-Modal Generation**

January 20, 2025

©2025 The Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)
209 N. Eighth Street
Stroudsburg, PA 18360
USA
Tel: +1-570-476-8006
Fax: +1-570-476-0860
acl@aclweb.org

ISBN 979-8-89176-213-8

Messages from the Organizers

Multimodal generation techniques have heralded new possibilities in creative content generation. Yet, the evaluation of such multimodal outputs remains a largely uncharted area, with fundamental questions still unresolved. These include understanding the contributions of individual modalities, the utility of pre-trained large language models in multimodal contexts, and the metrics for assessing faithfulness and fairness in generated outputs.

The first EvalMG workshop seeks to address these gaps by convening leading minds from natural language processing, computer vision, and multimodal AI. Our objective is to spearhead the development of robust evaluation methodologies that will propel further research in multimodal generation.

We received 21 submissions for this workshop, out of which 7 were accepted, including 5 long papers and 2 short papers. We have invited 11 reviewers and each submission was rigorously reviewed by at least three reviewers. The meta-reviews and final decisions were collaboratively handled by the organizing team.

We extend our deepest gratitude to all contributors—authors, reviewers, and particularly The University of Adelaide, for their generous support of this workshop. Your collective efforts are instrumental in shaping the future of multimodal research.

Organizing Committee

- Wei Emma Zhang
- Xiang Dai
- Desmond Elliot
- Byron Fang
- Haojie Zhuang
- Mong Yuan Sim
- Weitong Chen

- Necva Bölücü
- Danae Sanchez Villegas
- Wenhao Liang
- Lipin Guo
- Ali Shakeri
- Yutong Qu

Table of Contents

<i>A Dataset for Programming-based Instructional Video Classification and Question Answering</i> Sana Javaid Raja, Adeel Zafar and Aqsa Shoaib	1
<i>CVT5: Using Compressed Video Encoder and UMT5 for Dense Video Captioning</i> Mohammad Javad Pirhadi, Motahhare Mirzaei and Sauleh Eetemadi	10
<i>If I feel smart, I will do the right thing: Combining Complementary Multimodal Information in Visual Language Models</i> Yuyu Bai and Sandro Pezzelle	24
<i>LLaVA-RE: Binary Image-Text Relevancy Evaluation with Multimodal Large Language Model</i> Tao Sun, Oliver Liu, JinJin Li and Lan Ma	40
<i>Persian in a Court: Benchmarking VLMs In Persian Multi-Modal Tasks</i> Farhan Farsi, Shahriar Shariati Motlagh, Shayan Bali, Sadra Sabouri and Saeedeh Momtazi	52
<i>TaiwanVQA: A Benchmark for Visual Question Answering for Taiwanese Daily Life</i> Hsin-Yi Hsieh, Shang Wei Liu, Chang Chih Meng, Shuo-Yueh Lin, Chen Chien-Hua, Hung-Ju Lin, Hen-Hsen Huang and I-Chen Wu	57
<i>Guiding Vision-Language Model Selection for Visual Question-Answering Across Tasks, Domains, and Knowledge Types</i> Neelabh Sinha, Vinija Jain and Aman Chadha	76

