

第2回講演会のお知らせ

Video Generation and Understanding with Multimodal Learning

Videos are inherently multimodal, blending static appearance, dynamic motion, physical geometry, and semantic context. Yet, traditional models often treat them as simple sequences of flat RGB frames. This talk explores how explicitly leveraging these underlying modalities can fundamentally advance both video generation and understanding. The speaker will introduce a series of frameworks designed to bridge these signals, including an appearance-motion diffusion distillation method that speeds up video sampling, a human-centric model that pairs appearance with depth, and a text-to-motion approach for synthesizing diverse human actions. Shifting to video understanding, the talk will highlight a debiased action recognition strategy that isolates true motion from misleading semantic cues to significantly improve model robustness. Together, these approaches demonstrate that integrating appearance, motion, geometry, and semantics is key to building more coherent and generalizable video models.

講師

Meta
Research Scientist

Yuanhao Zhai 様

Bio: Yuanhao Zhai is a Research Scientist at Meta, where he focuses on image and video generative modeling for creative and personalized advertising. He earned his Ph.D. in Computer Science and Engineering from the University at Buffalo, co-advised by Prof. David Doermann and Prof. Junsong Yuan. During his doctoral studies, Yuanhao investigated video understanding and video generation, publishing extensively in top-tier computer vision conferences and journals. Prior to his PhD, he received his B.Eng. in Computer Science from Xi'an Jiaotong University in 2020, where he was advised by Ziyi Liu and Prof. Le Wang at the Institute of Artificial Intelligence and Robotics.

2026

6/25 土

10:00-11:30

オンライン

参加申込

下記Googleフォームよりお申込みください。

<https://forms.gle/rb1P3UCrHPpCdhMa7>

