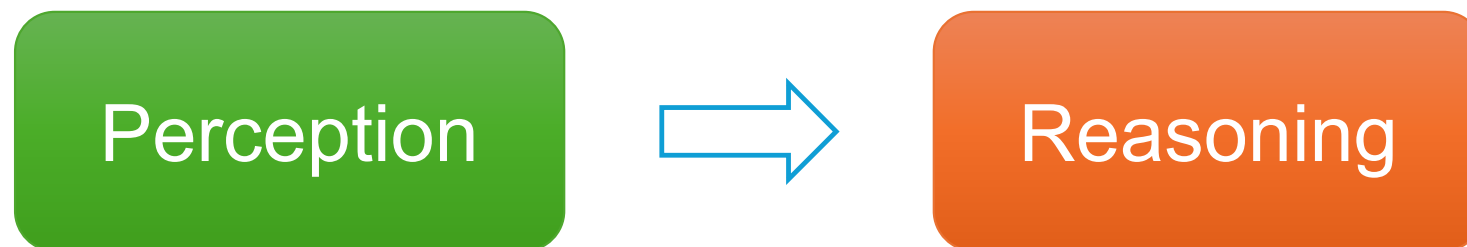ICCV 2025 Tutorial on

# Towards Comprehensive Reasoning in Vision-Language Models

Oct 19 8.30am- 12pm

# From Perception to Reasoning for VLM

Perception ⟹ Reasoning

- Pattern recognition ≠ Reasoning

- Reasoning means compositional structure, causal understanding, and multi-step inference.

- Compared with language-only models, VLMs still lag when visual evidence must be integrated step-by-step.

# Bring Comprehensive Reasoning Capabilities into VLM

Pattern → Process Reasoning: Modeling reasoning as a process rather than a shortcut.

Single shot → Structured and Reliable Reasoning: Pursuing consistency and robustness for trustworthy inference.

Static → Interactive and Agentic Reasoning: Grounding reasoning in real-world perception and action.

# Our Recent Explorations on Process Reasoning



In **FrameMind**, we studied how interleaved visual-textual chains to enforce multi-step compositional reasoning via Reinforcement Learning for VLM video understanding.

Ge, H., Wang, Y., Chang, K. W., Wu, H., & Cai, Y. (2025). FameMind: Frame-Interleaved Video Reasoning via Reinforcement Learning. arXiv e-prints, arXiv-2509.

# Our Recent Explorations on Structured and Reliable Reasoning



MRFD for VLM hallucinations [EMNLP 2025]

CHAINMPQ for VLM Relation Understanding [arxiv 2025]

We proposes works like **MRFD** and **CHAINMPQ** to enhance reliability and self-consistency to make reasoning more trustworthy.

# Our Recent Explorations on Interactive and Agentic Reasoning



Through **Dimo-GUI** [EMNLP2025] and **Vistawise** [EMNLP2025], we explore how reasoning emerges when a model interacts with its environment.

# Tutorial Schedule

| Time | Session | Speaker |
|---|---|---|
| 8:30 - 8:35 | Opening Remark: Motivation and Overview [Abstract] [Slides] | Yujun Cai |
| 8:35 - 9:10 | Invited Talk: LMMs-Lab: Building Multimodal Intelligence [Abstract] [Slides] | Ziwei Liu |
| 9:10 - 9:35 | Video-TT Challenge: Towards Advanced Video Reasoning and Understanding [Abstract] [Slides] | Yuhao Dong, Yuanhan Zhang, Ziwei Liu, and Representative Teams |
| 9:35 - 10:10 | Invited Talk: Reasoning in Multimodal GUI Agents: An Exploration-Driven Perspective [Abstract] [Slides] | Chi Zhang |
| 10:10 - 10:45 | Invited Talk: Mathematical Reasoning in Visual Contexts [Abstract] [Slides] | Kai-Wei Chang |
| 10:45 - 11:20 | Invited Talk: Chain-of-Look Visual Reasoning [Abstract] [Slides] | Junsong Yuan |
| 11:20 - 11:55 | Invited Talk: Grounding Anything in Images and Videos for Comprehensive Reasoning [Abstract] [Slides] | Ming-Hsuan Yang |
| 11:55 - 12:00 | Closing Remark [Abstract] [Slides] | Yiwei Wang |

# Organizers & Speakers

Yujun Cai
@UQ

Jun Liu
@Lancaster

Yiwei Wang
@UCM

Ming-Hsuan Yang
@UCM

Kai-Wei Chang
@UCLA

Junsong Yuan
@UB

Ziwei Liu
@NTU

Chi Zhang
@Westlake

Yuanhan Zhang
@NTU

Yuhao Dong
@NTU