



AAAI 2025 Tutorial T04
Time: 2025-02-25 8:30-12:30
Location: Room 118A

Foundation Models Meet Embodied Agents



Manling Li
Northwestern



Yunzhu Li
Columbia



Jiayuan Mao
MIT



Wenlong Huang
Stanford



Northwestern
University



COLUMBIA



Stanford
University



AAAI 2025 Tutorial T04
Time: 2025-02-25 8:30-12:30
Location: Room 118A

Remaining Challenges

AAAI Tutorial: Foundation Models Meet Embodied Agents



Northwestern
University



COLUMBIA



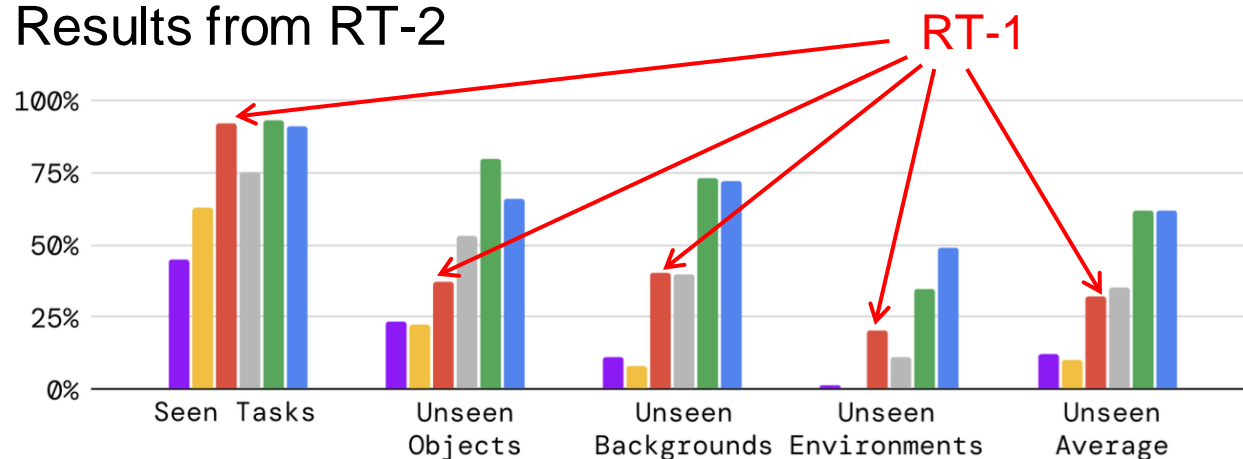
Stanford
University

- ❑ Inconsistent results across papers

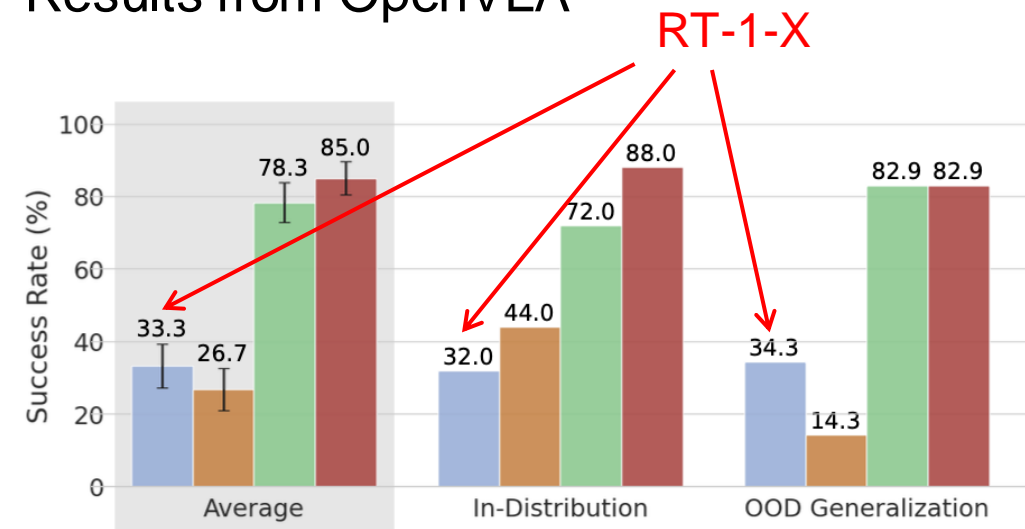
Results from RT-1

Model	Seen Tasks	Unseen Tasks	Distractors	Backgrounds
Gato (Reed et al., 2022)	65	52	43	35
BC-Z (Jang et al., 2021)	72	19	47	41
BC-Z XL	56	43	23	35
RT-1 (ours)	97	76	83	59

Results from RT-2



Results from OpenVLA



- Evaluation is primarily conducted in the real world
 - Real-world evaluation is costly and noisy
 - “We have large enough budget such that we can still make progress.”
 - Weak correlation between training loss and real-world success rate.
 - Training objectives vs task-specific metrics, training vs testing horizons



ALOHA 2

- ❑ What about evaluation in simulation?
 - ❑ Sim-to-real gap: rigid / deformable / cloth
 - ❑ Efficient asset generation
 - ❑ Digitalization of the real world
 - ❑ Procedural generation of realistic and diverse scenes
 - ❑ Correlation between sim and real

ImageNet in Embodied AI?

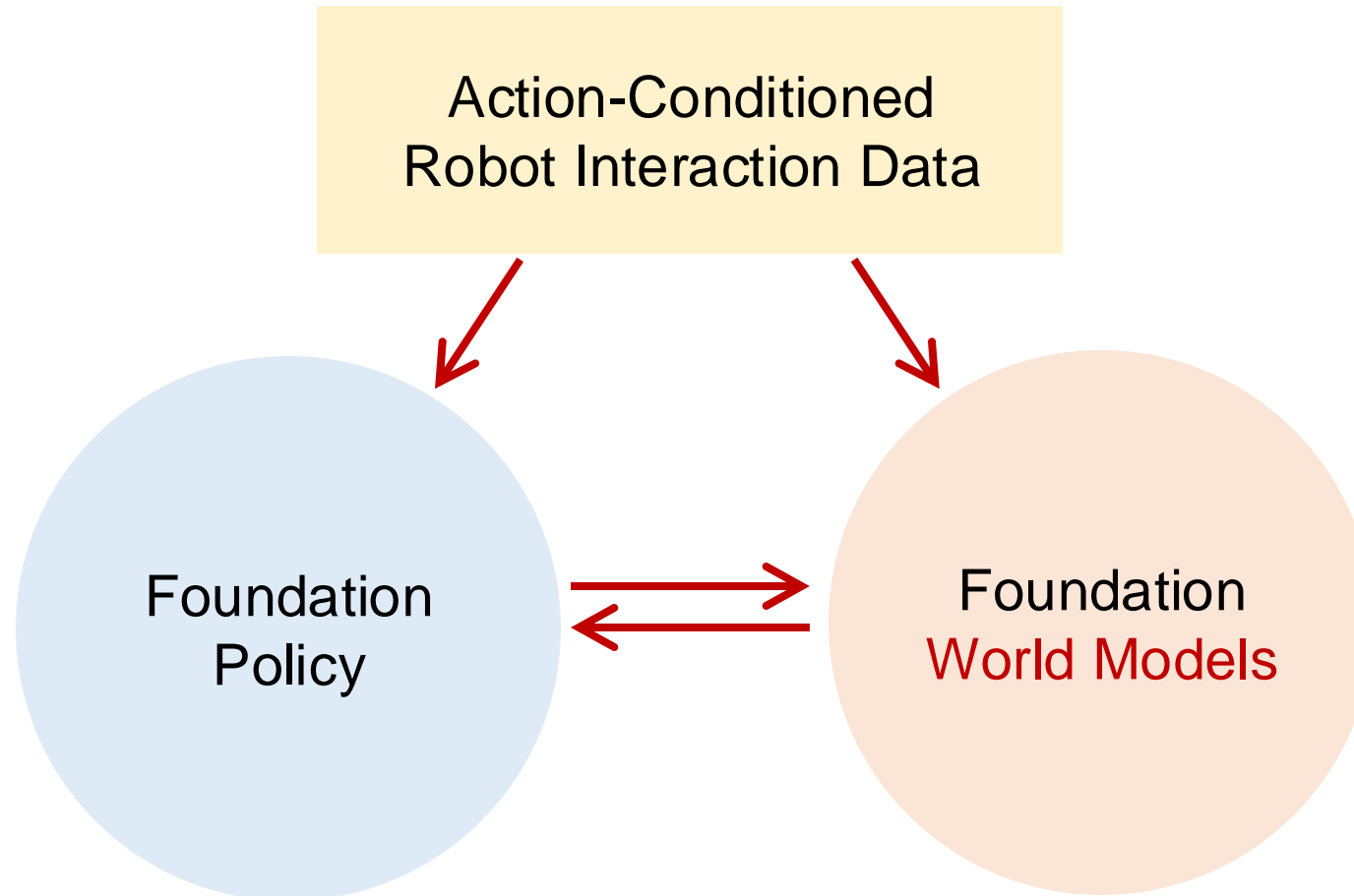


BEHAVIOR-1K

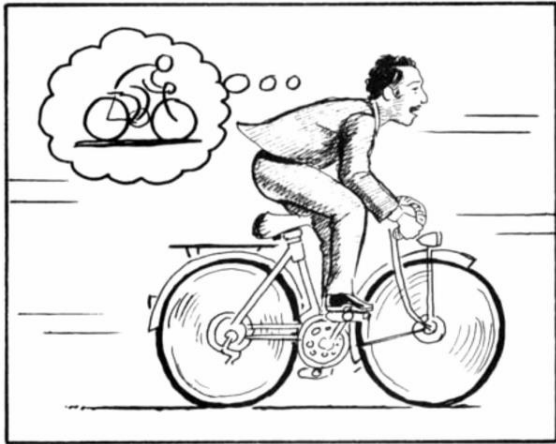


Habitat 3.0

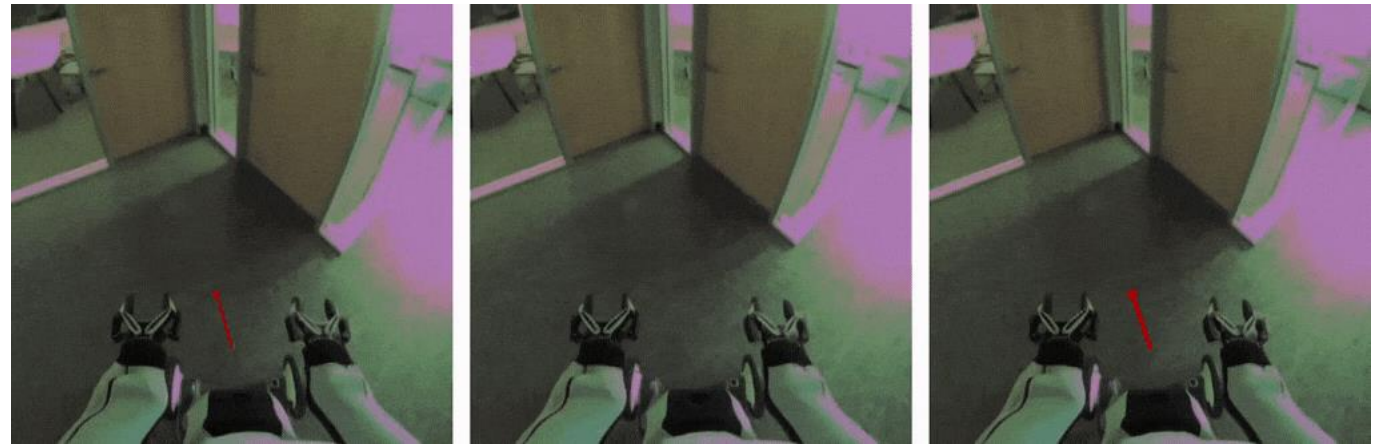
- My definition of world models: **action-conditioned future prediction**



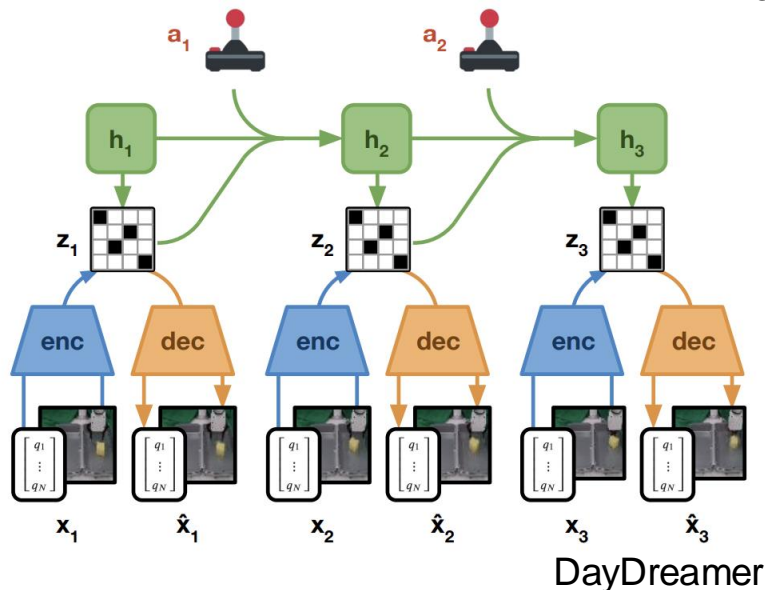
- My definition of world models: **action-conditioned future prediction**



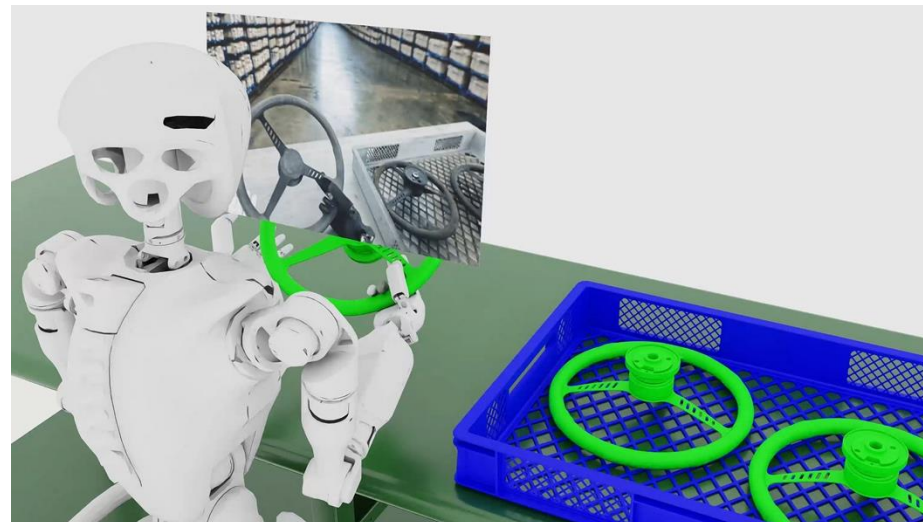
World Models



1X World Models



DayDreamer



Nvidia Cosmos - World Foundation Model

- 3D?
- Structural Prior?
- Learning + Physics?
- Corr. w/ Real World

- ❑ Current foundation models are not tailored for embodied agents
 - ❑ LLM/VLM can fail in embodied-related tasks
 - ❑ Limited understanding of geometric / embodied / physical interactions
 - ❑ Reinforcement learning (RL) from human feedback → RL from [Embodied Feedback](#)



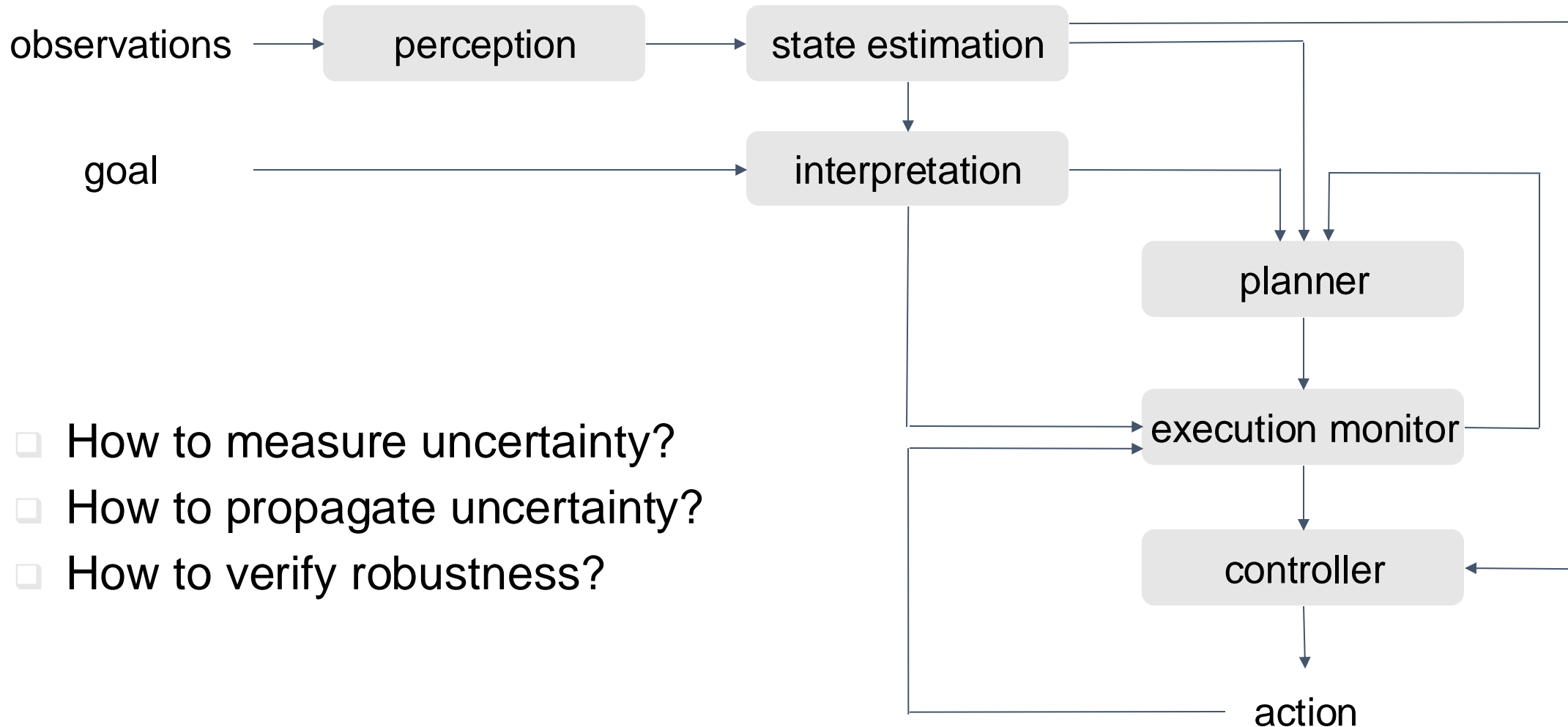
GPT



Segment Anything



DINOv2

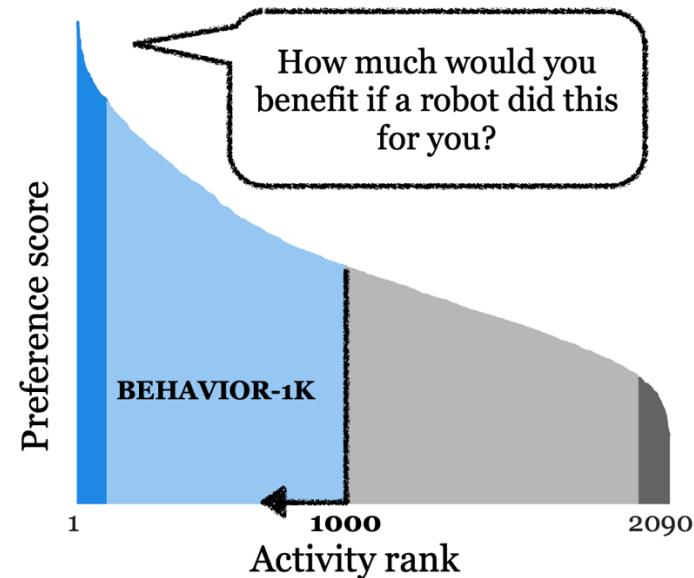


- ❑ How to measure uncertainty?
- ❑ How to propagate uncertainty?
- ❑ How to verify robustness?

- ❑ Adapt to new scenarios
- ❑ Adapt to human preferences
- ❑ Self improve / life-long learning



Adapt to new scenarios

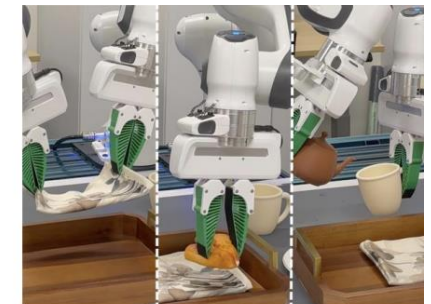
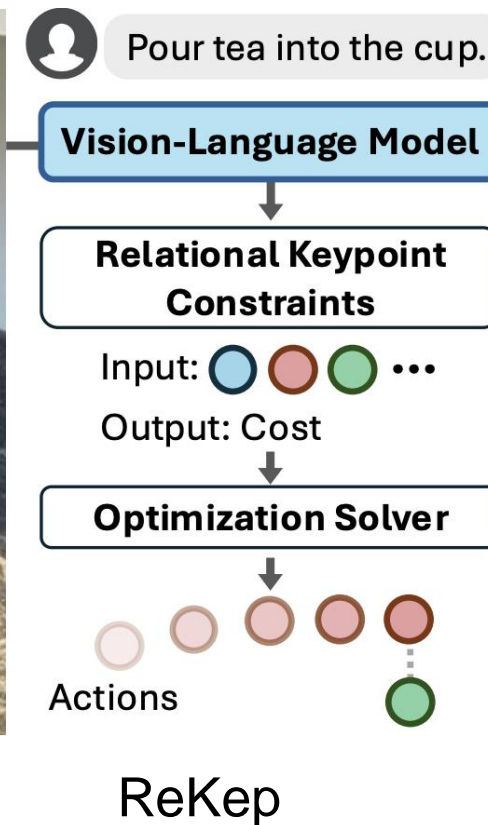
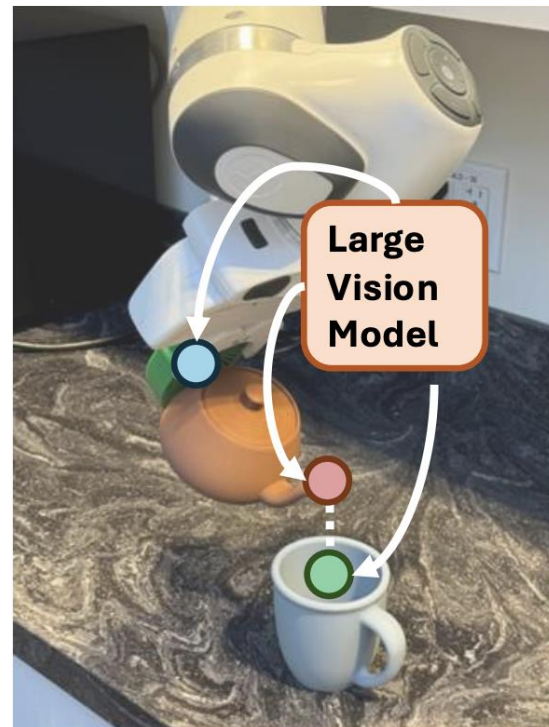
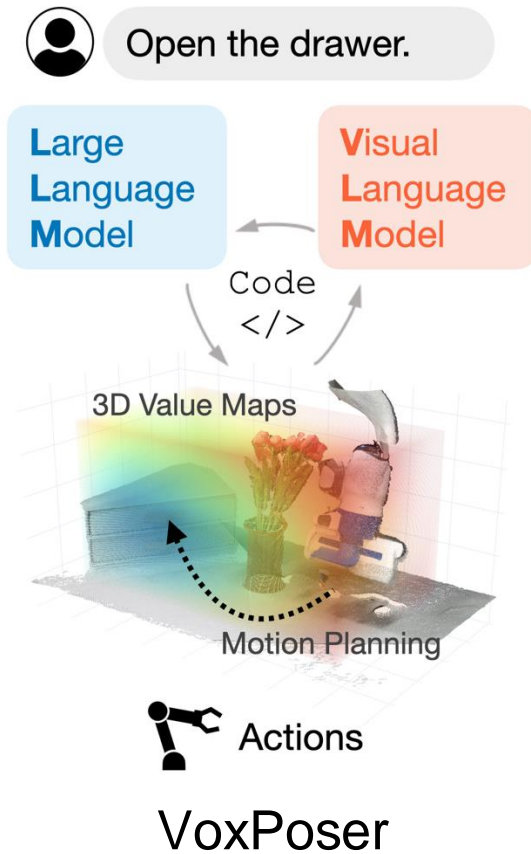


Adapt to human preferences



Improve through experience

- ❑ Every robotics work is a system work
- ❑ System-level considerations:
 - ❑ delays / computing / modules talking to each other (high-level vs low-level)



Multi-Stage



Bimanual

- ❑ Every robotics work is a system work
- ❑ System-level considerations: delays / computing / modules talking to each other

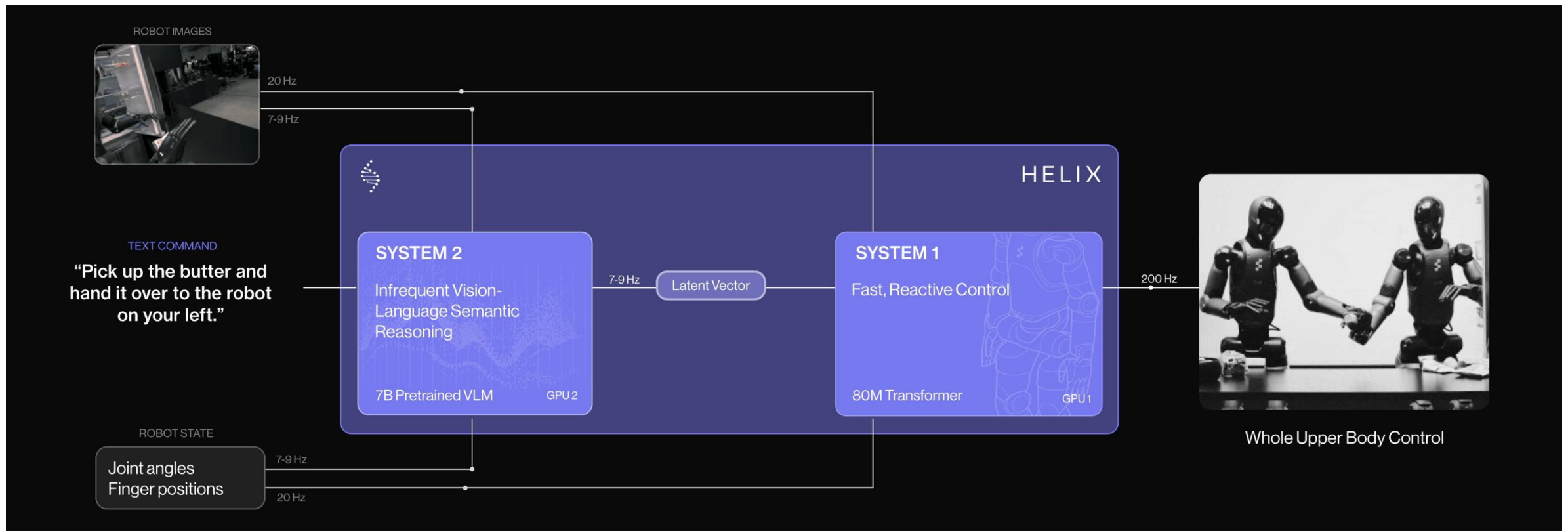


Figure A1: Helix

- ❑ Bring together researchers working on both virtual & physical agents
- ❑ Approaching the problem from a more structured lens: MDP
- ❑ Key techniques, emerging opportunities, and notable challenges

The background of the slide is a dark, blue-toned image featuring a futuristic, metallic robot on the right side. The robot has a human-like form with visible joints and a head. The background is filled with abstract, glowing digital patterns and data-like elements, creating a high-tech, artificial intelligence atmosphere.

Foundation Models Meet Embodied Agents

@ AAI 2025 Tutorial

08:30 - 12:30, Feb 25, 2025

118 A, Pennsylvania Convention Center, Philadelphia, Pennsylvania



AAAI 2025 Tutorial T04
Time: 2025-02-25 8:30-12:30
Location: Room 118A

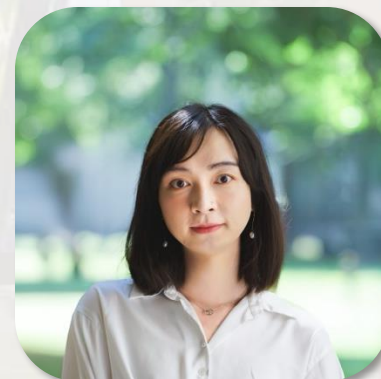
Foundation Models Meet Embodied Agents



Manling Li
Northwestern



Yunzhu Li
Columbia



Jiayuan Mao
MIT



Wenlong Huang
Stanford



Northwestern
University



COLUMBIA



Stanford
University