



# Building General Embodied Agents in the Physical World

Baoxiong Jia  
BIGAI

Figure generated by Nano Banana Pro

# About me

[buzz-beater.github.io](https://buzz-beater.github.io)



**Peking University**  
**B.S. in CS**  
**2014-2018**



**UCLA**  
**Ph.D. in CS**  
**2018-2022**



**BIGAI**  
**Research Scientist**  
**2022-Present**





@XRoboHub

# Embodied AI

The embodiment hypothesis is the idea that intelligence emerges in the interaction of an agent with an environment and as a result of sensorimotor activity

Smith & Gasser, The Development of Embodied Cognition: Six Lessons from Babies, 2005

Hardware, control, and locomotion



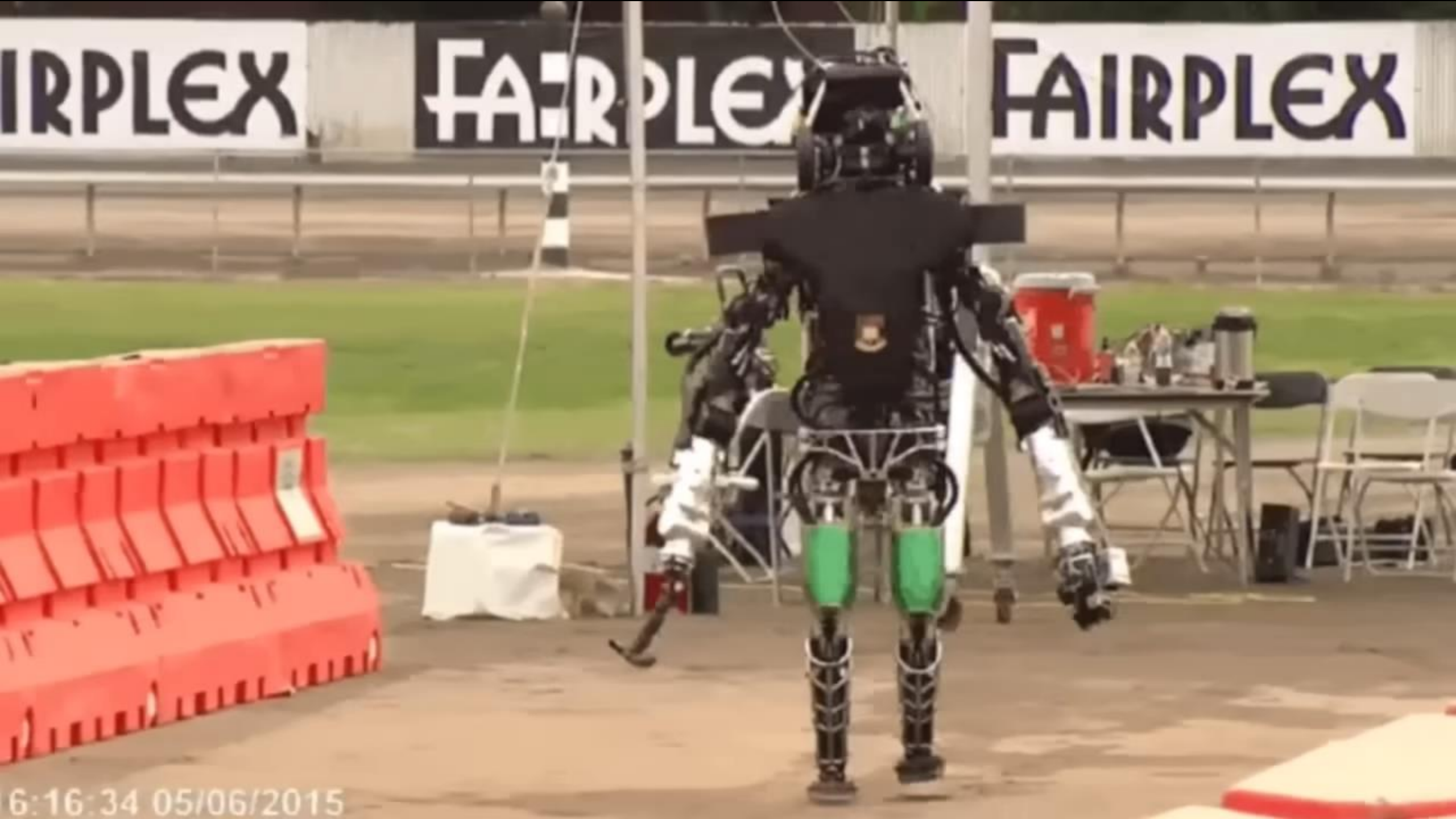
Boston Dynamics, Atlas | Partners in Parkour, 2022  
<https://www.youtube.com/watch?v=tF4DML7FIWk>

Interaction, reason, and plan



Damen et al., Scaling Egocentric Vision: The Epic-Kitchens Dataset, 2018



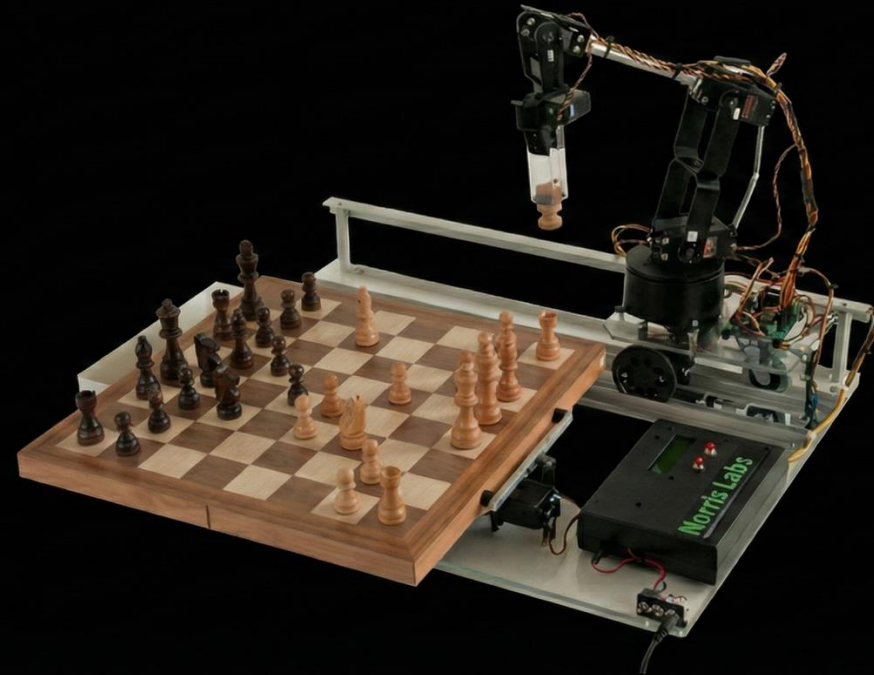
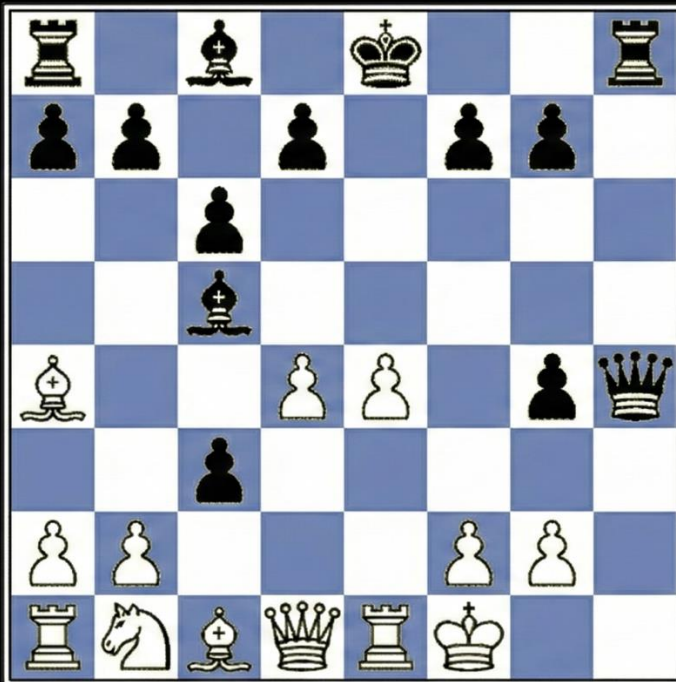


6:16:34 05/06/2015

# Moravec's Paradox

It's comparatively easy to make computers exhibit adult level performance on intelligence tests or playing checkers, and difficult or impossible to give them skills of a one-year-old when it comes to perception and mobility.

Hans Moravec, Mind Children, Harvard University Press 1988

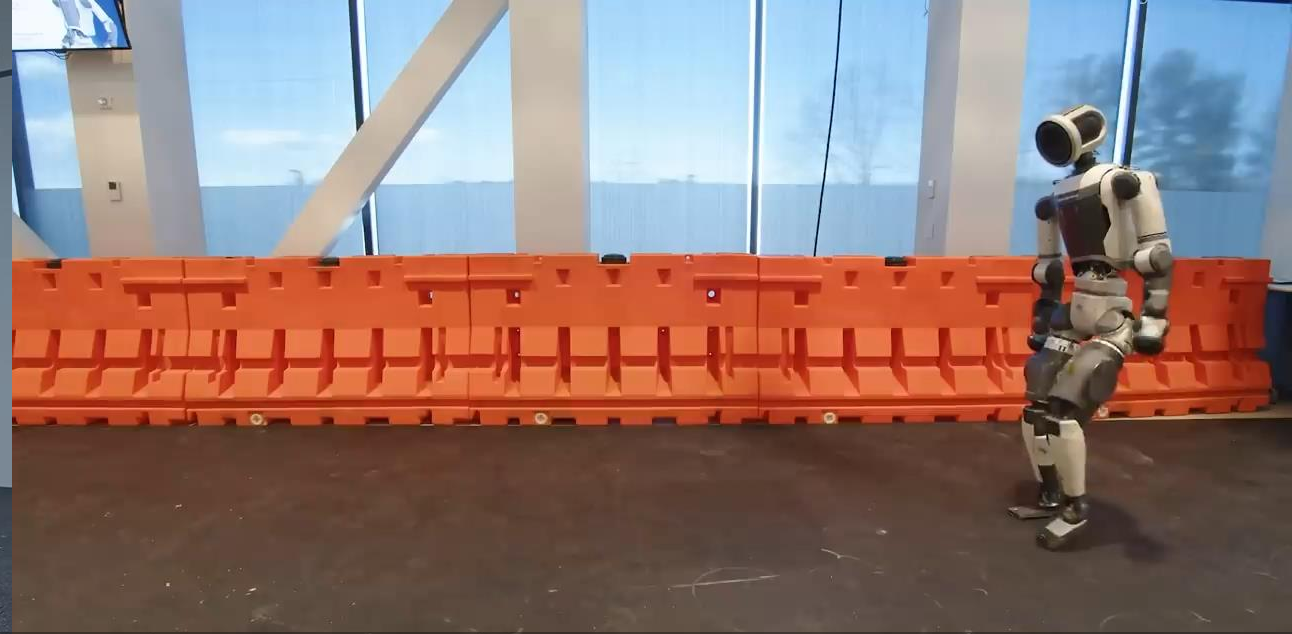




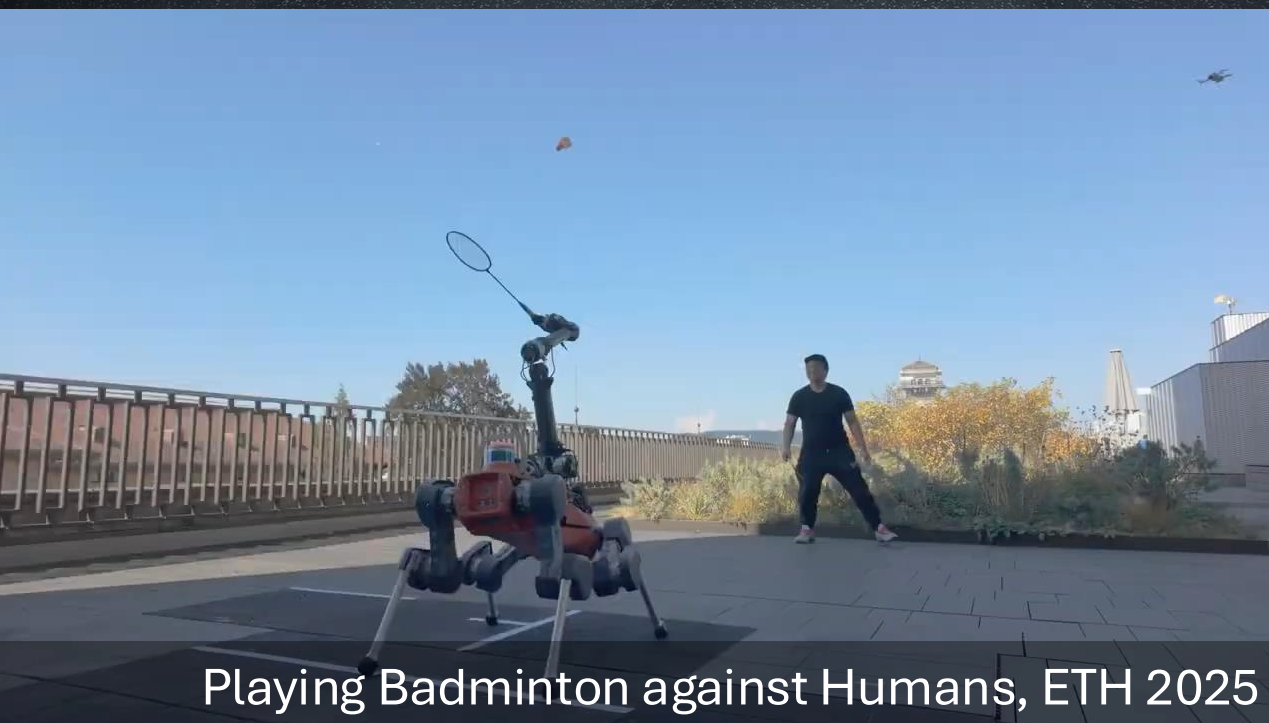
No speed-up in this video



UniTree Kungfu Kid 6.0, UniTree 2025



Walk, Run, Crawl, RL Fun, Boston Dynamics 2024



Playing Badminton against Humans, ETH 2025

Double spin into  
a 3.5 spin handstand



UniTree B2-W Talent Awakening, UniTree 2024



GPT/Qwen

1.2B Hours

$\pi_0$

10k Hours

autonomous 1

GR00T

88 Hours



# EMBODIED AI

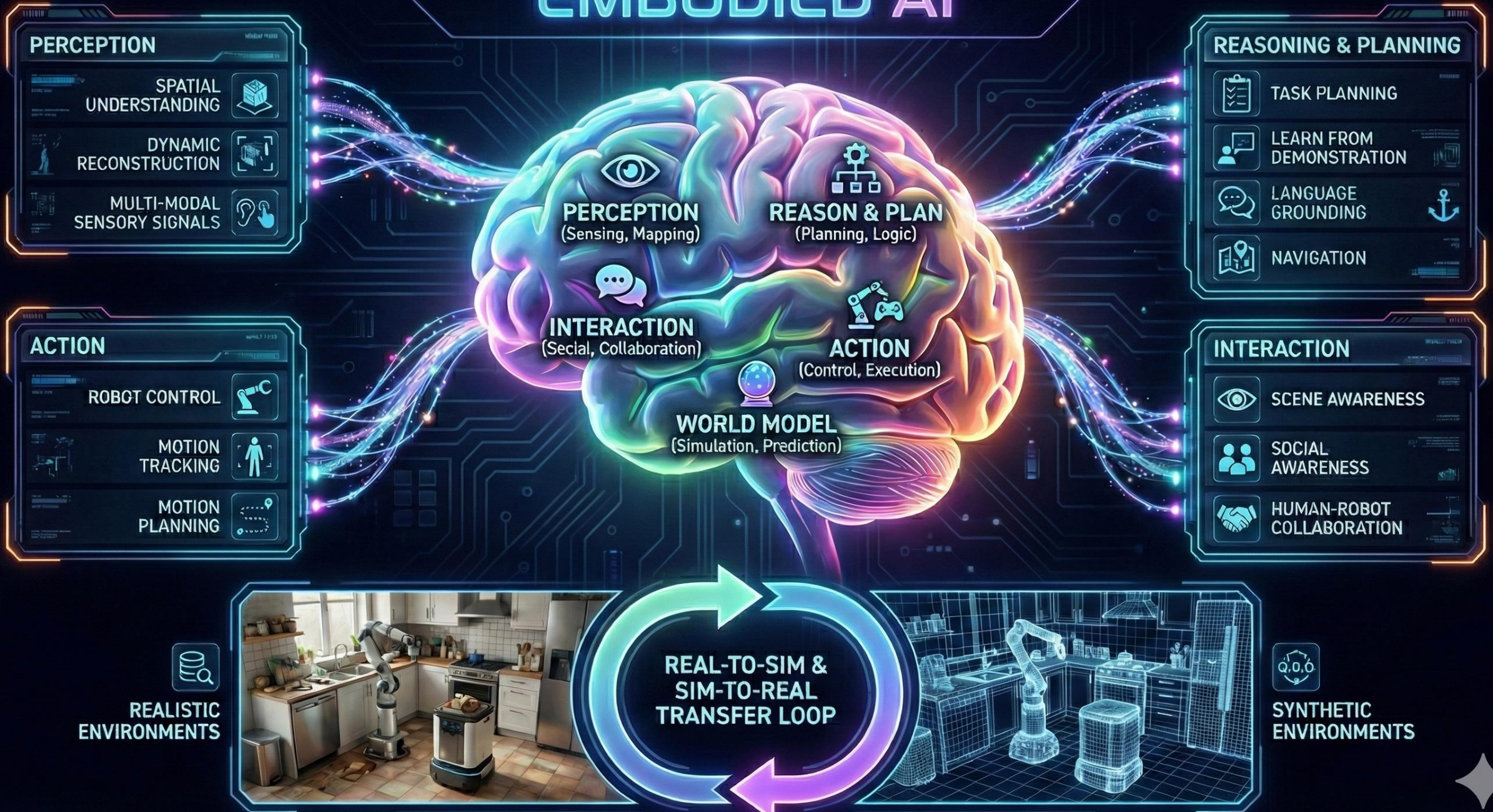


Figure generated by Nano Banana Pro



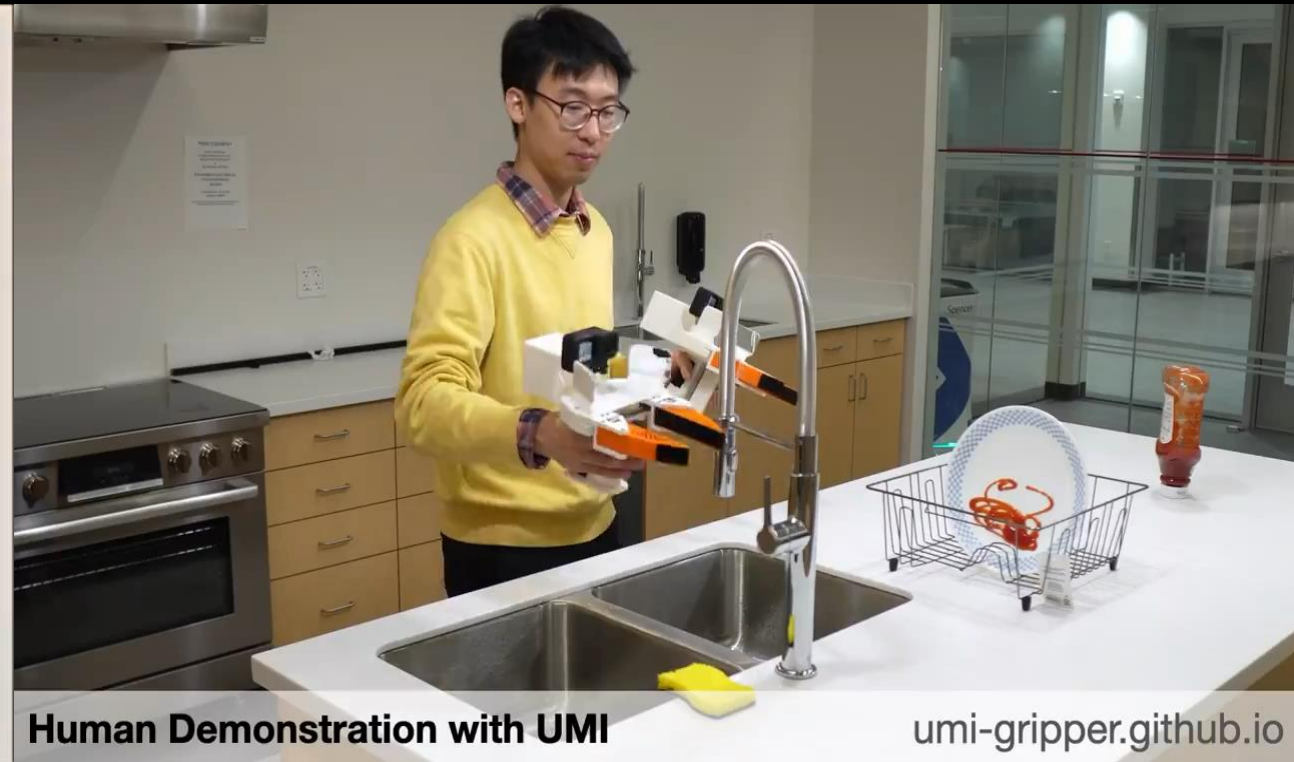
# 1. Building Interactable 3D Scenes for Embodied AI



# Collecting Data and Training Robot in the Real World is Expensive



**Damage the environment and objects**



**Low-cost hardware emerging but still low efficiency**

# Goal of Environment Creation

- High-quality **appearance** understanding for grounding and reasoning
- Fine-grained **geometry** understanding for simulation and physics
- Solid **dynamics** understanding for interaction and planning



eractable Replicas

Dynamic Reconstruction

culated Objects

an Splatting

**ArtGS (ICLR'25)**

**Digital Cousin Creation**

**MetaScenes (CVPR'25)**

**3D Scene Generation**

**SceneWeaver (NeurIPS'25)**

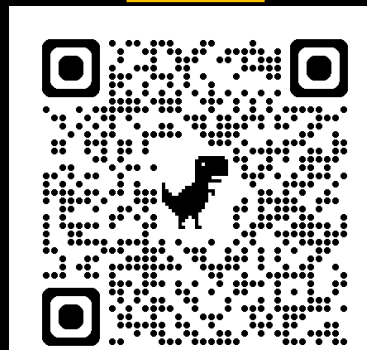
**Best Paper, RoboGen@IROS25**



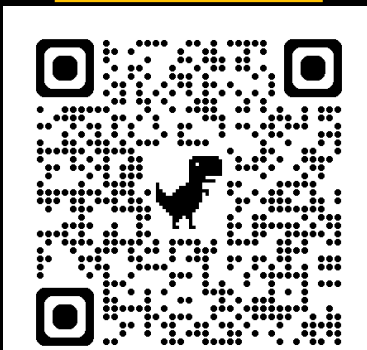
# 4D Gaussians for Dynamic Reconstruction

- *(ICLR'25) Building Interactable Replicas of Complex Articulated Objects via Gaussian Splatting*
- *(ArXiv'25) VideoArtGS: Building Digital Twins of Articulated Objects from Monocular Video*

ArtGS

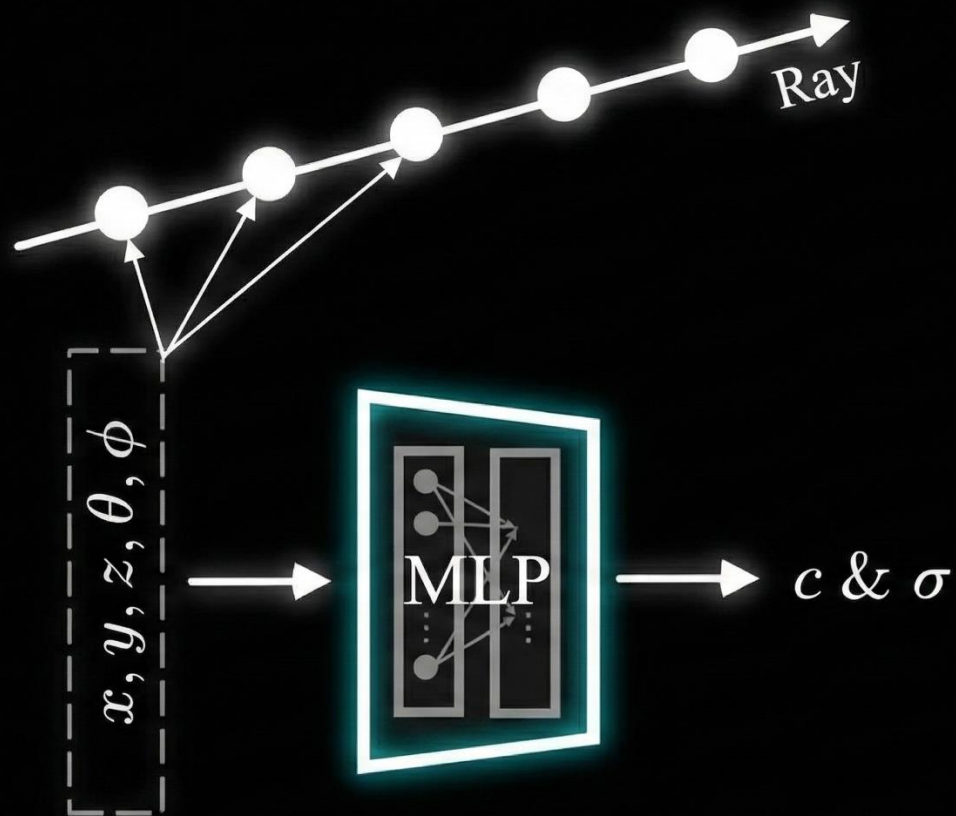


VideoArtGS

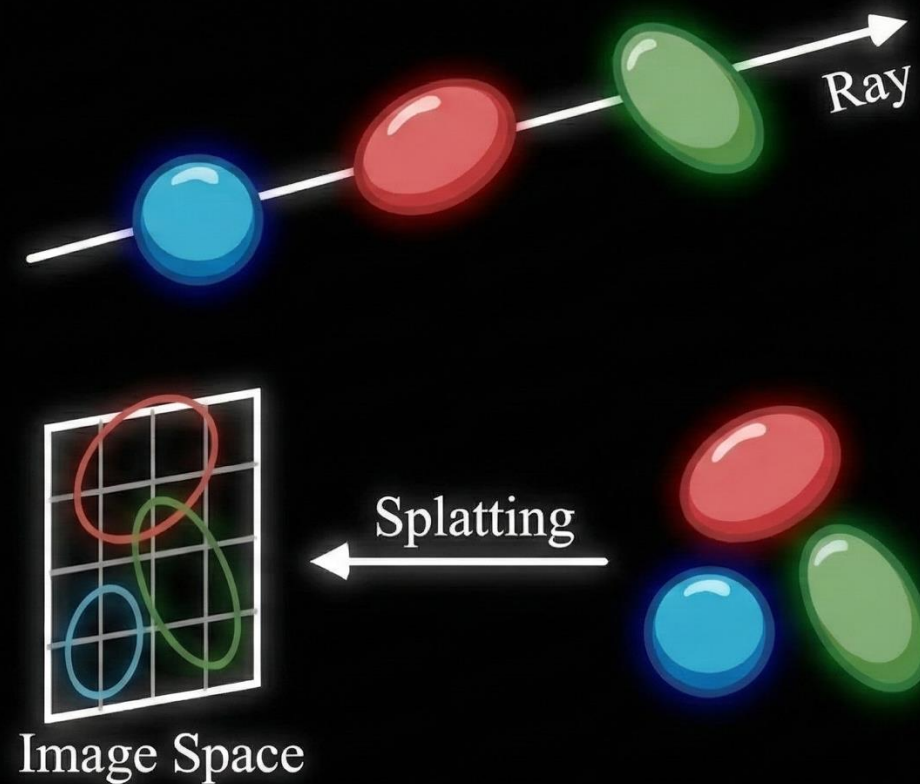




# Static 3D Scene Reconstruction



(a) NeRF



(b) 3D GS

**Learn 3D by projecting to multiple views as supervision**

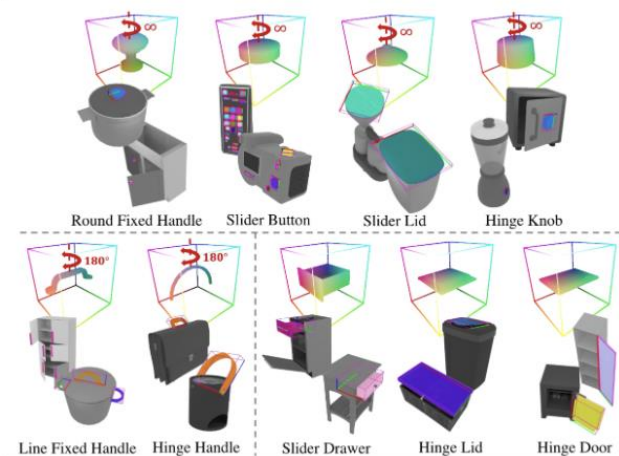
# Manipulation Involves Dynamic Objects



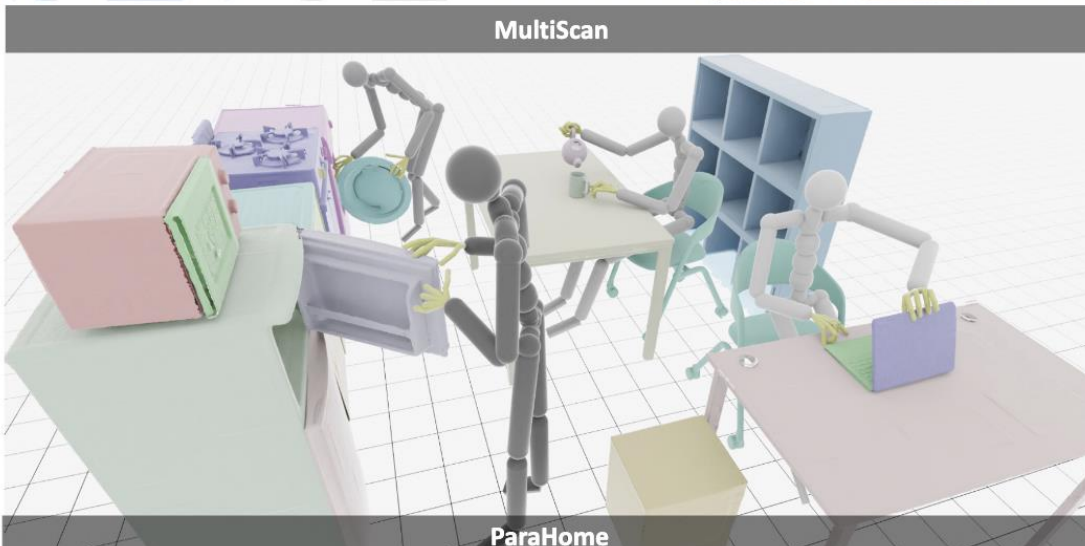
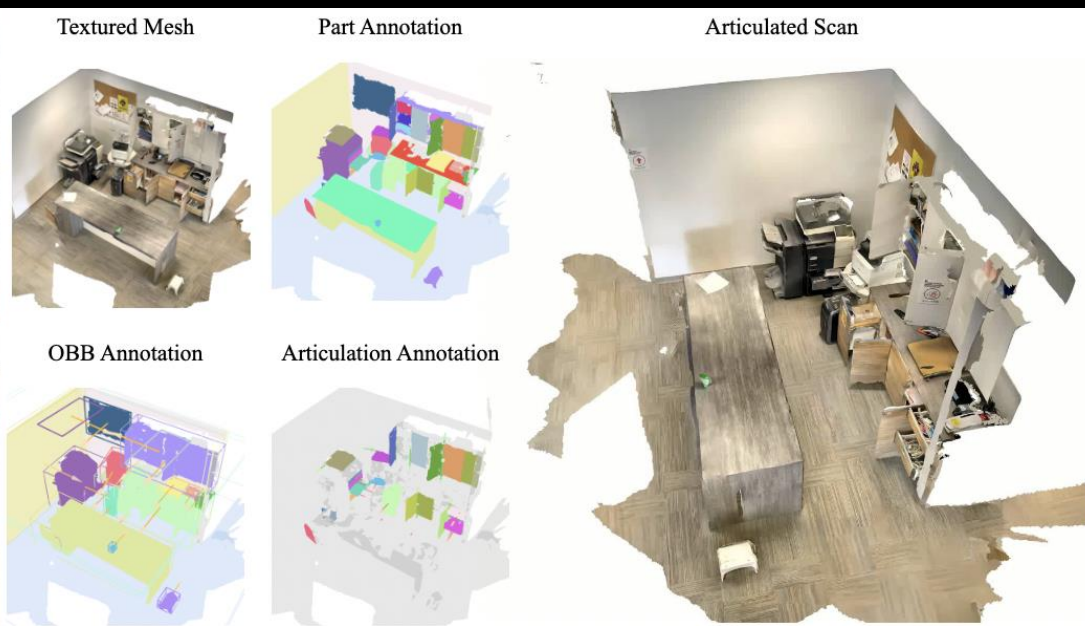
PartNet-Mobility



AKB-48



GAPartNet



ParaHome

In reality, we deal with dynamic, **articulated objects** whose **geometry and shape change** during interaction, making them difficult to reconstruct



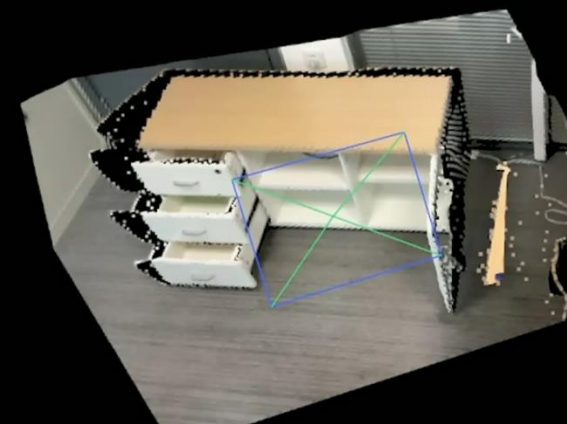
# Image Supervision is Ambiguous for Articulation Learning

**Key Challenge:** The observed pixel motion results from four entangled factors:

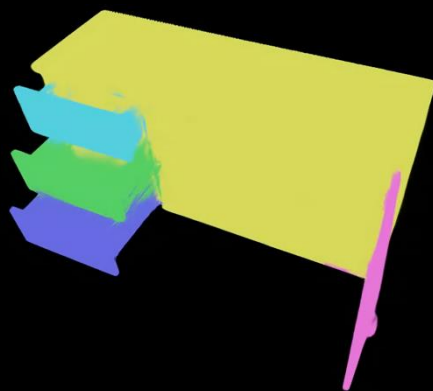
Camera  
trajectory



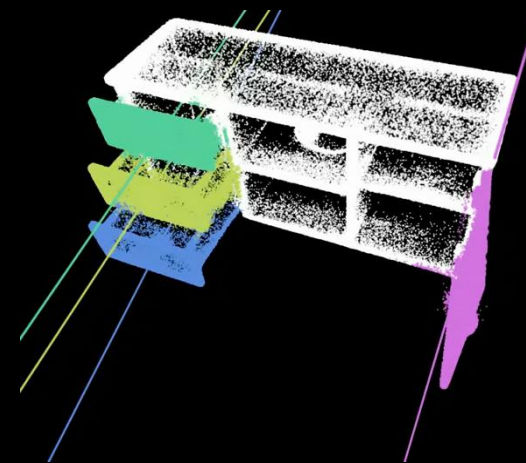
Object  
Geometry



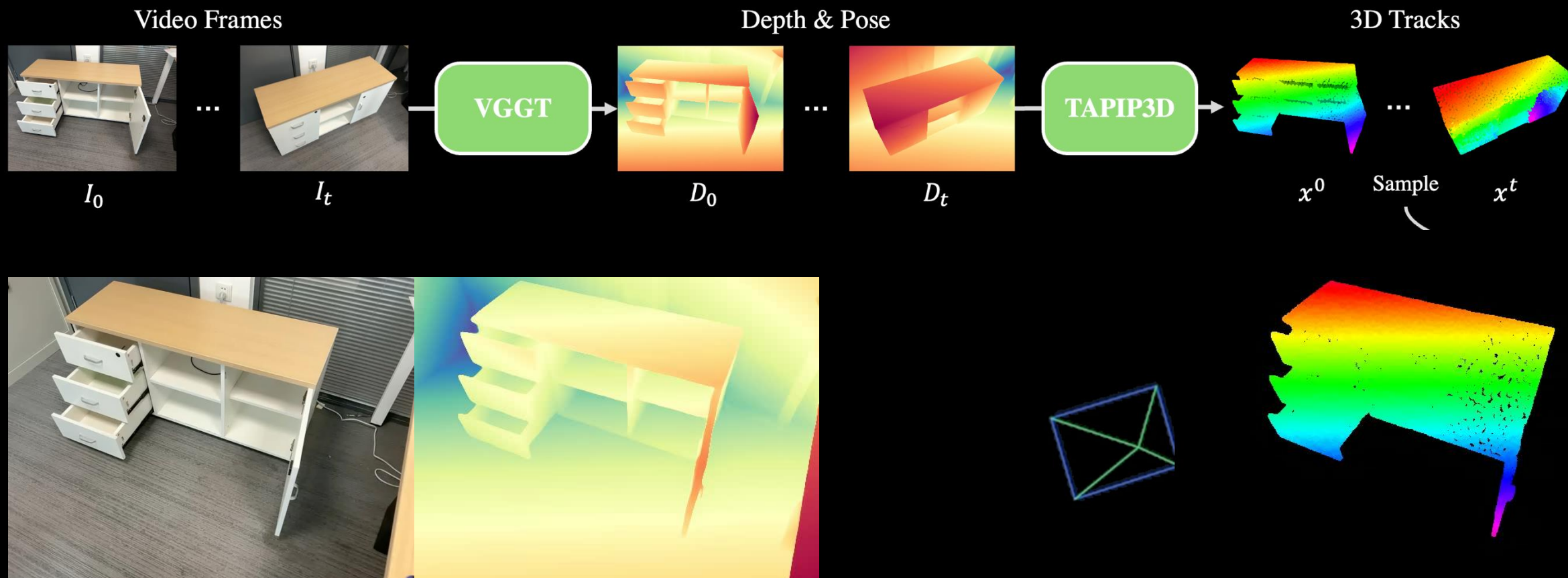
Part  
Segmentation



Articulation  
Dynamics



# Camera, Depth, Tracks Estimation





# Geometry Reconstruction & Articulation Learning

Video Frames



$I_0$

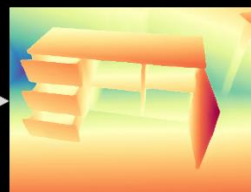
...



$I_t$



VGGT



$D_0$

Depth & Pose

...

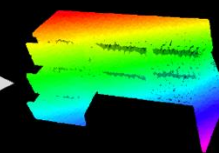


$D_t$



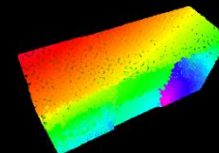
TAPIP3D

3D Tracks

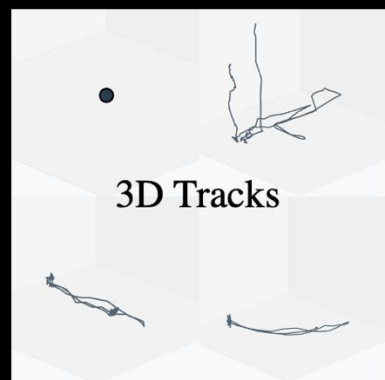


$x^0$

Sample

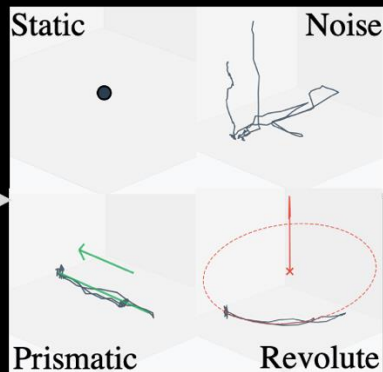


$x^t$



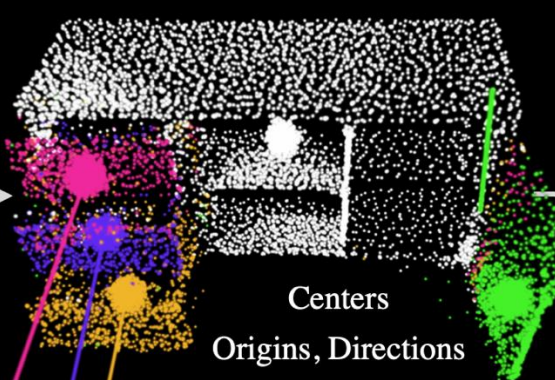
3D Tracks

Motion  
Analysis



Static Noise  
Prismatic Revolute

Motion  
Clustering



Centers  
Origins, Directions

Init



Deformation  
Field

$$\hat{x}^{t_2} = \mathcal{F}(\hat{x}^c, t_2)$$

$$\mathcal{L} = (\hat{x}^{t_2} - x^{t_2})^2$$



First  
N-frames

Init



Canonical Gaussians

$\mathcal{G}^c$



Articulation-based  
Deformation

$\mathcal{G}^t$



Deformed Gaussians

Render



$$\hat{I}_t \longleftrightarrow I_t$$



$$\hat{D}_t \longleftrightarrow D_t$$



# Real-world Experiments



Cabinet

Laptop

Cabinet

Microwave

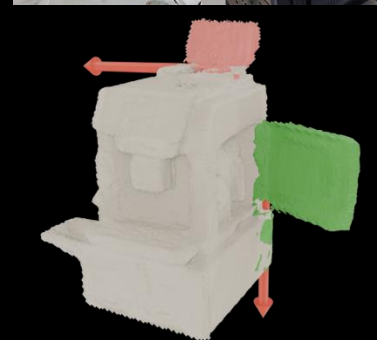
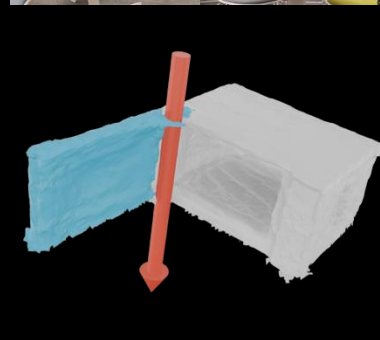
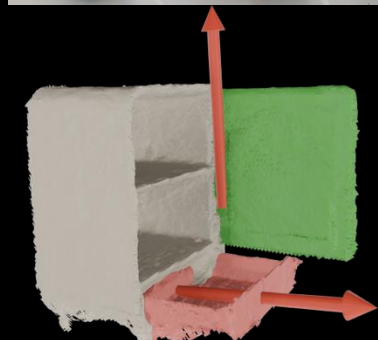
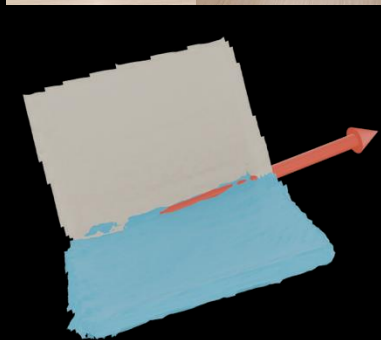
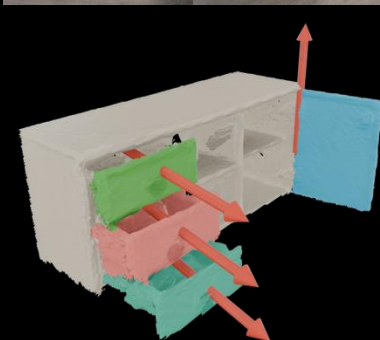
Coffee Machine

Chair

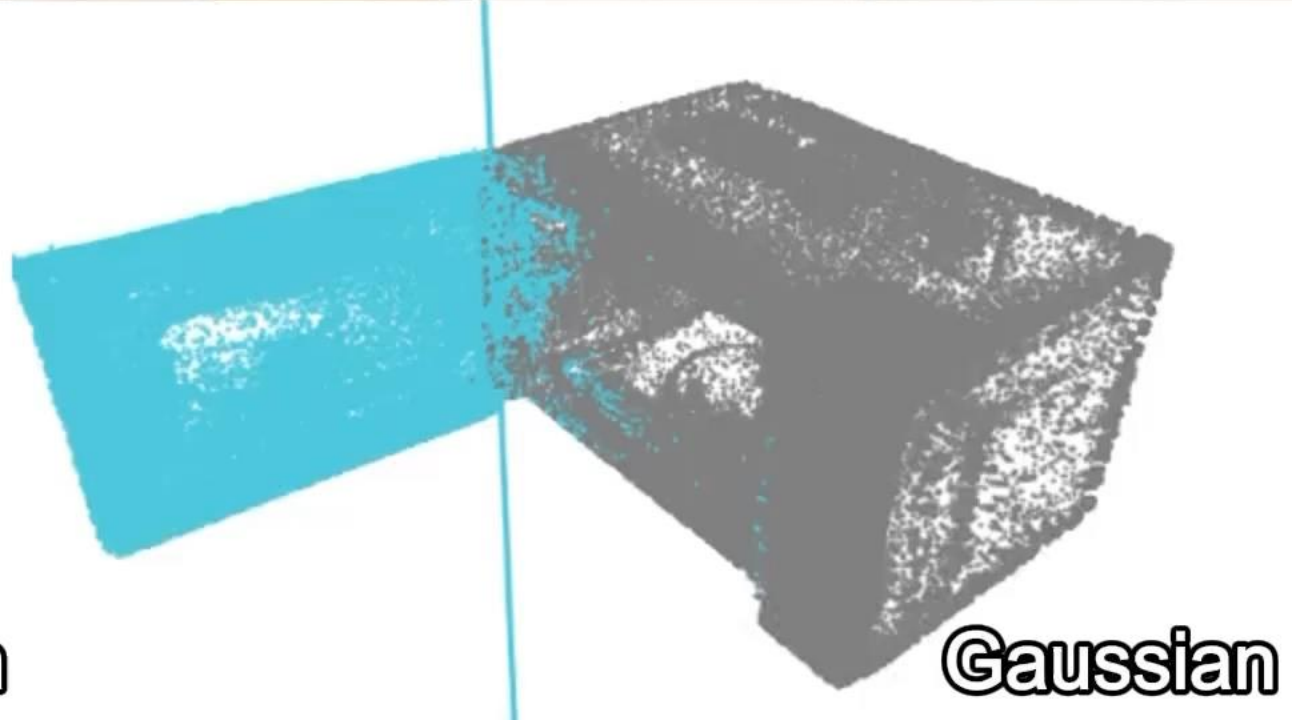
Input  
Frames



Results







# Limitations & Takeaways

- **Difficult to scale-up**
  - ❖ Feed-forward reconstruction
  - ❖ Active camera trajectory selection
- **Limited quality for simulation**
  - ❖ **Physical priors** (e.g. plane) during reconstruction
  - ❖ Existing assets and **generative models** as guidance for reconstruction
  - ❖ **System vs. Model?**

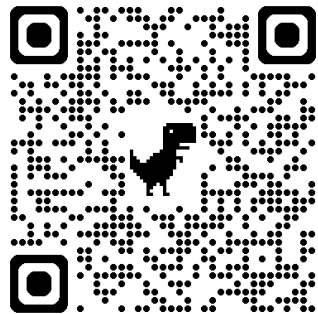
ArtGS: Building Interactable Replicas  
of Complex Articulated Objects  
via Gaussian Splatting



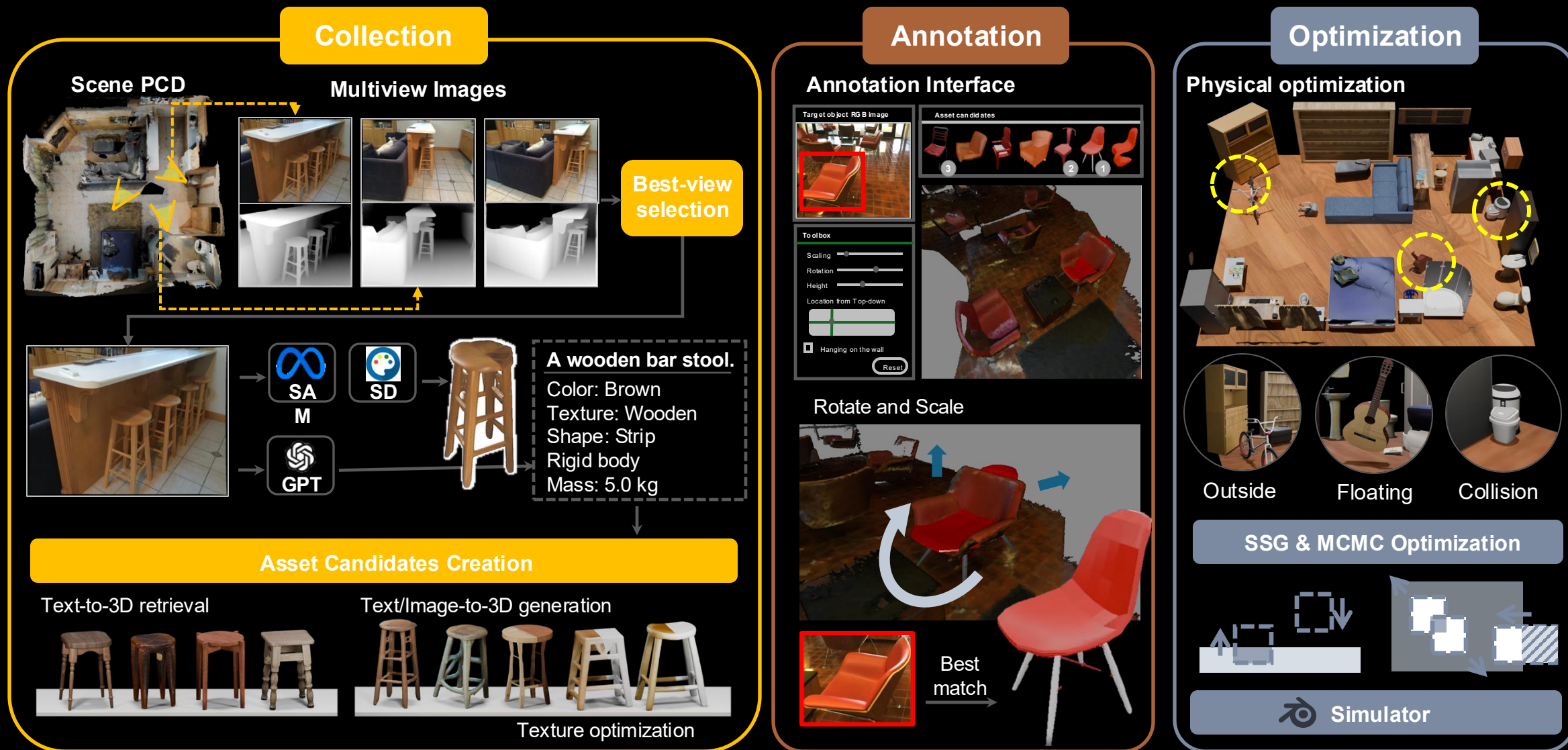
# Digital Cousin Creation with 3D AIGC

- *(CVPR'25) MetaScenes: Towards Automated Replica Creation for Real-World 3D Scans*

[MetaScenes](#)



# Leveraging AIGC & Online Assets for Digital Twin Creation

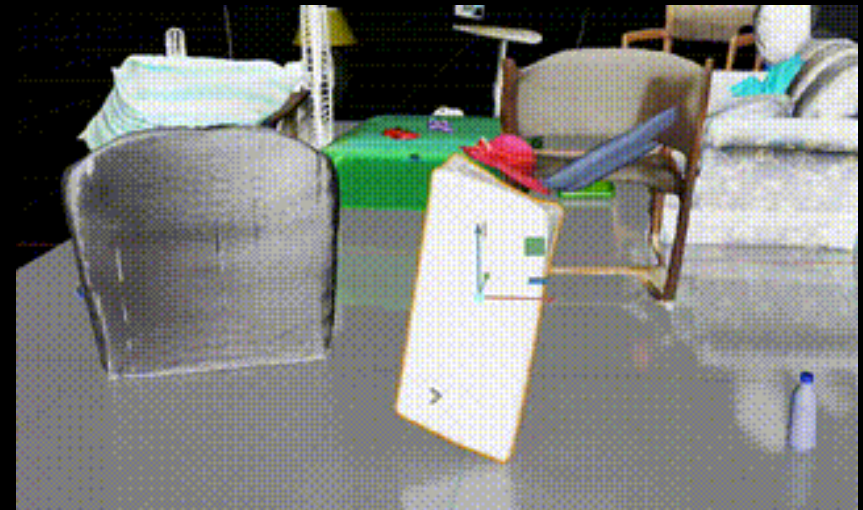
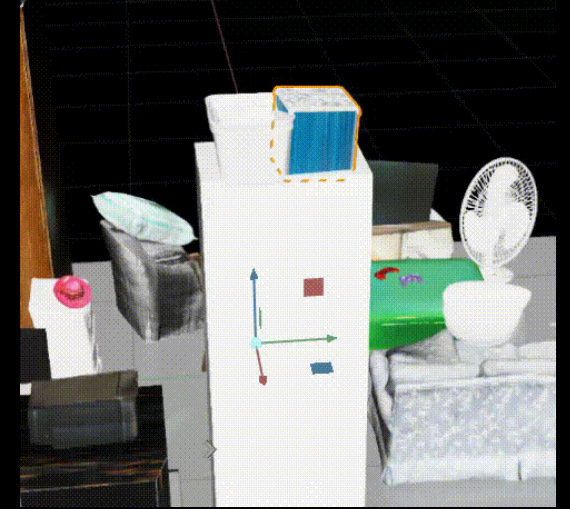






# Limitations & Takeaways

- **Unstable generation quality**
  - ❖ **Physical-based** post-optimization
  - ❖ Better 3D generative models (e.g. SAM3D)
- **Sequential error** in the pipeline
  - ❖ Better orchestration of tools
  - ❖ **Iterative refinement** of the generated results?





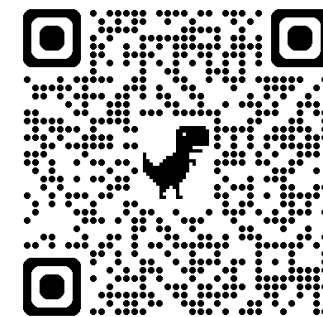


# Agentic Tool-Use for 3D Scene Generation

- (NeurIPS'25) SceneWeaver: All-in-One 3D synthesis with an Extensible and Self-Reflective Agent (*Best Paper, RoboGen@IROS'25*)



SceneWeaver



## Tool Cards

Initializer



Tool #1: MetaScenes



Tool #2: PhyScene

Room Type  
Room Size  
3D Layout  
Relation  
Assets

Implementer



Tool #1: ACDC



Tool #2: LLM

Refiner



Tool #1: VLM  
Update Rotation



Tool #2: Rule  
Add Relation



Tool #3: LLM  
Remove Object

Query



STOP



## Planner

Think

Reflection

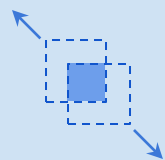
Action

Feedback

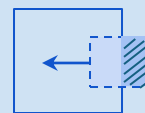


## Executor

+Physical optimize



Collision



Out of bound

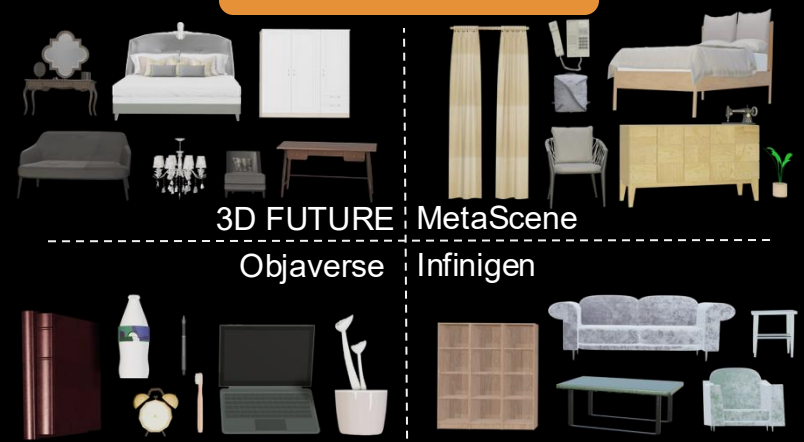


Relation

## Final Scene

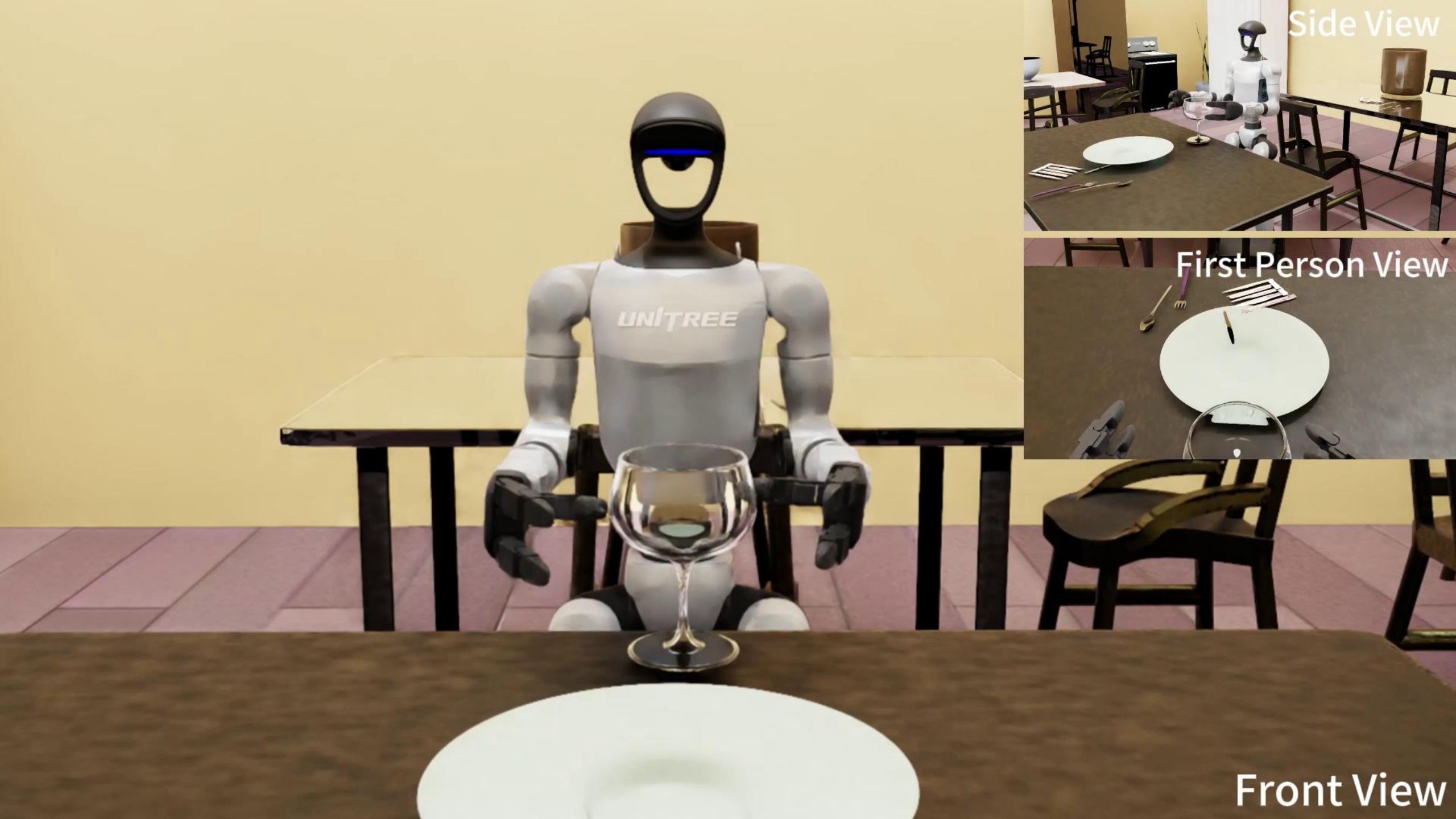


## Assets









Side View

First Person View

Front View



## 2. Bridging Reasoning and Action for General Embodied AI Agents

REASONING CAPTIONING:  
"ASSEMBLING GEAR  
MECHANISM"

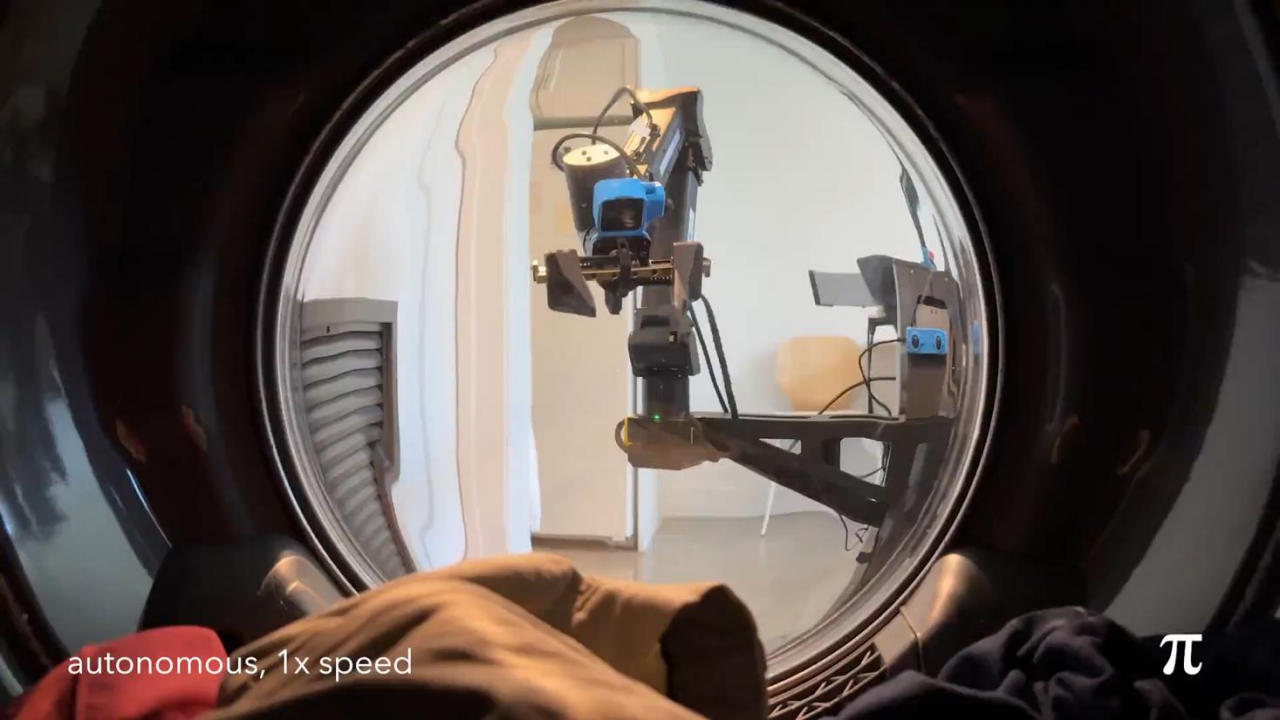
SCENE  
UNDERSTANDING

SEMANTIC  
SEGMENTATION

OBJECT DETECTION  
& MANIPULATION

Figure generated by Nano Banana Pro



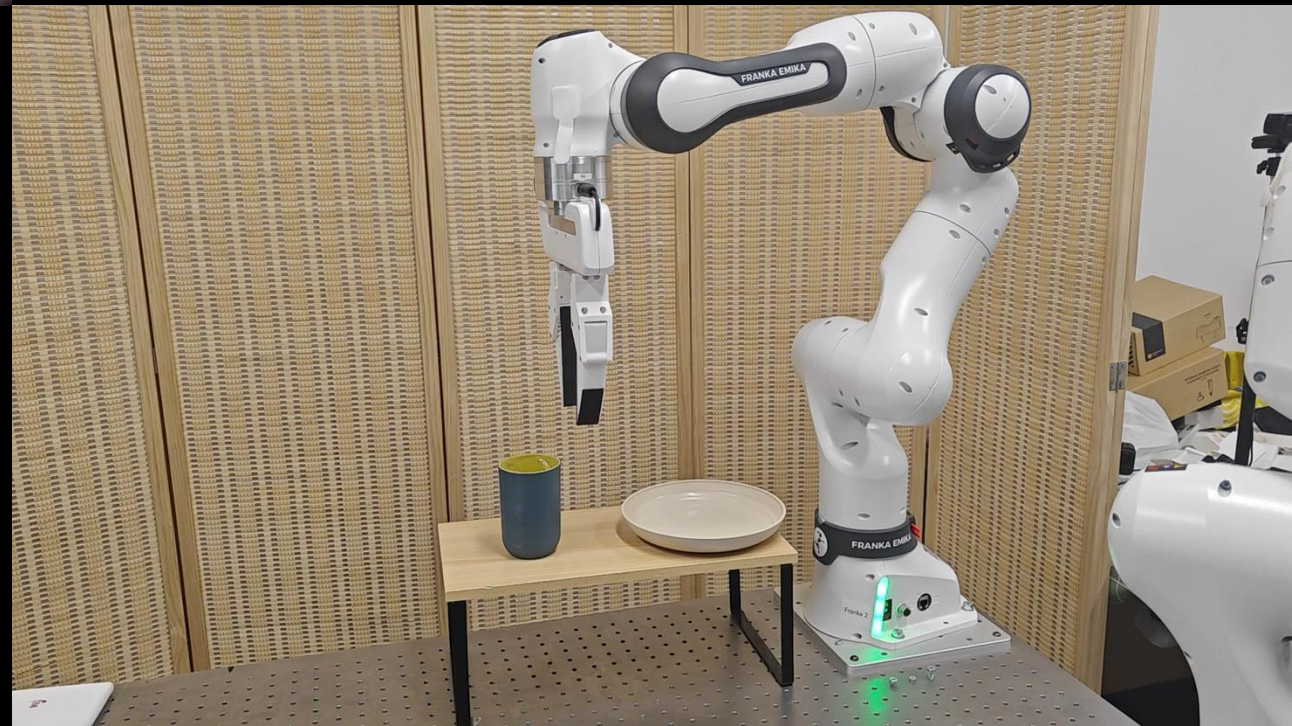


## Training Generalist Policies

- Leveraging large-scale pre-trained VLMs
- Pre-trained with large-scale data
- Still limited generalizability on tasks and embodiments

## Adapting to Your Specific Scene

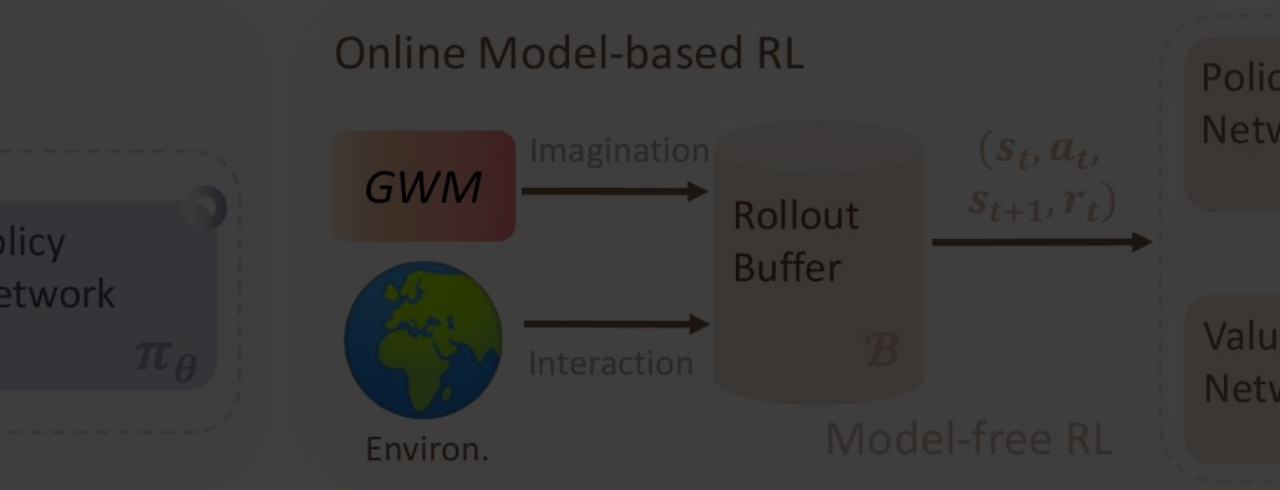
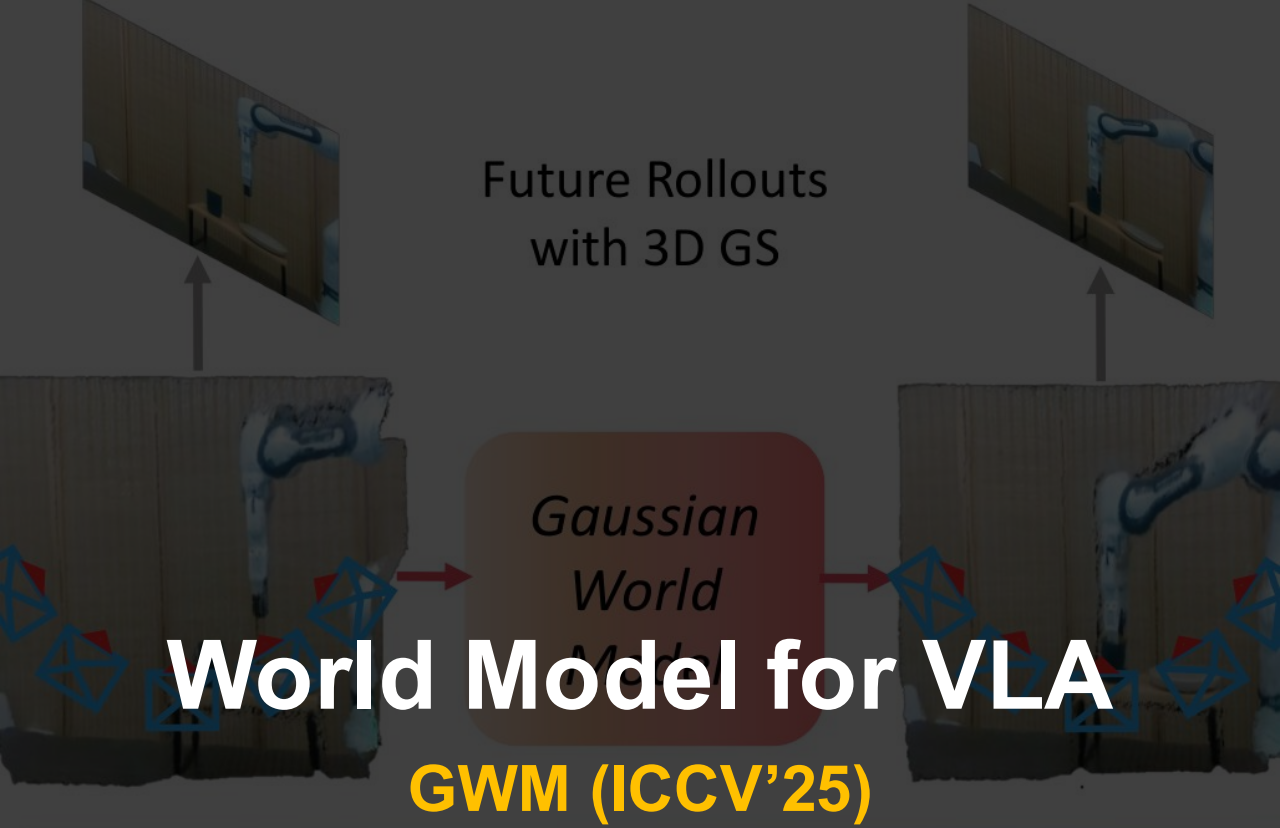
- Rolling out “almost” successful trajectories but hard to improve
- Can only afford few-shot demonstrations
- Sensitive to capturing modalities and viewpoints





# Goal of Reasoning and Planning

- Capable of finishing **diverse tasks** according to instructions
- Enable **spatial understanding** for existing VLMs for better backbones
- **Efficient representation** for effective learning from high-cost data



**An Embodied Generalist Agent in 3D World**

Jiangyong Huang<sup>1,2\*</sup>, Silong Yong<sup>1,3\*</sup>, Xiaojian Ma<sup>1\*</sup>, Xiongkun Linghu<sup>1\*</sup>, Puhao Li<sup>1,4</sup>, Yan Wang<sup>1</sup>, Qing Li<sup>1</sup>, Song-Chun Zhu<sup>1,2,4</sup>, Baoxiong Jia<sup>1</sup>, Siyuan Huang<sup>1</sup>

<sup>1</sup>Beijing Institute for General Artificial Intelligence (BIGAI)  
<sup>2</sup>Peking University <sup>3</sup>Carnegie Mellon University <sup>4</sup>Tsinghua University

<https://embodied-generalist.github.io/>

**Embodied General Agent**

Capabilities: Perception, Grounding, Reasoning, Planning, Acting

**LEO (ICML'24)**

3D Object Captioning, 3D Question Answering, Dialogue Task Planning, Embodied Reasoning, Embodied Navigation, Robotic Manipulation

**3D World**

**LEO**

The block contains information about the "Embodied General Agent LEO". It lists the authors and their affiliations, a QR code, and a website link. The agent's capabilities are listed as Perception, Grounding, Reasoning, Planning, and Acting. Below this, specific tasks are listed: 3D Object Captioning, 3D Question Answering, Dialogue Task Planning, Embodied Reasoning, Embodied Navigation, and Robotic Manipulation. A cartoon lion character holding a sign that says "LEO" is featured. At the bottom, there are several 3D world visualizations and a small image of a robot arm.

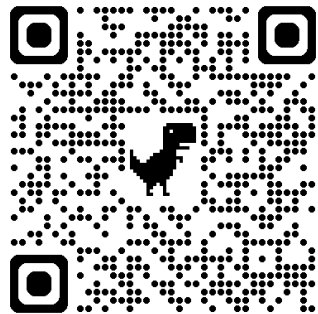


# Open-Vocabulary Mobile Manipulation

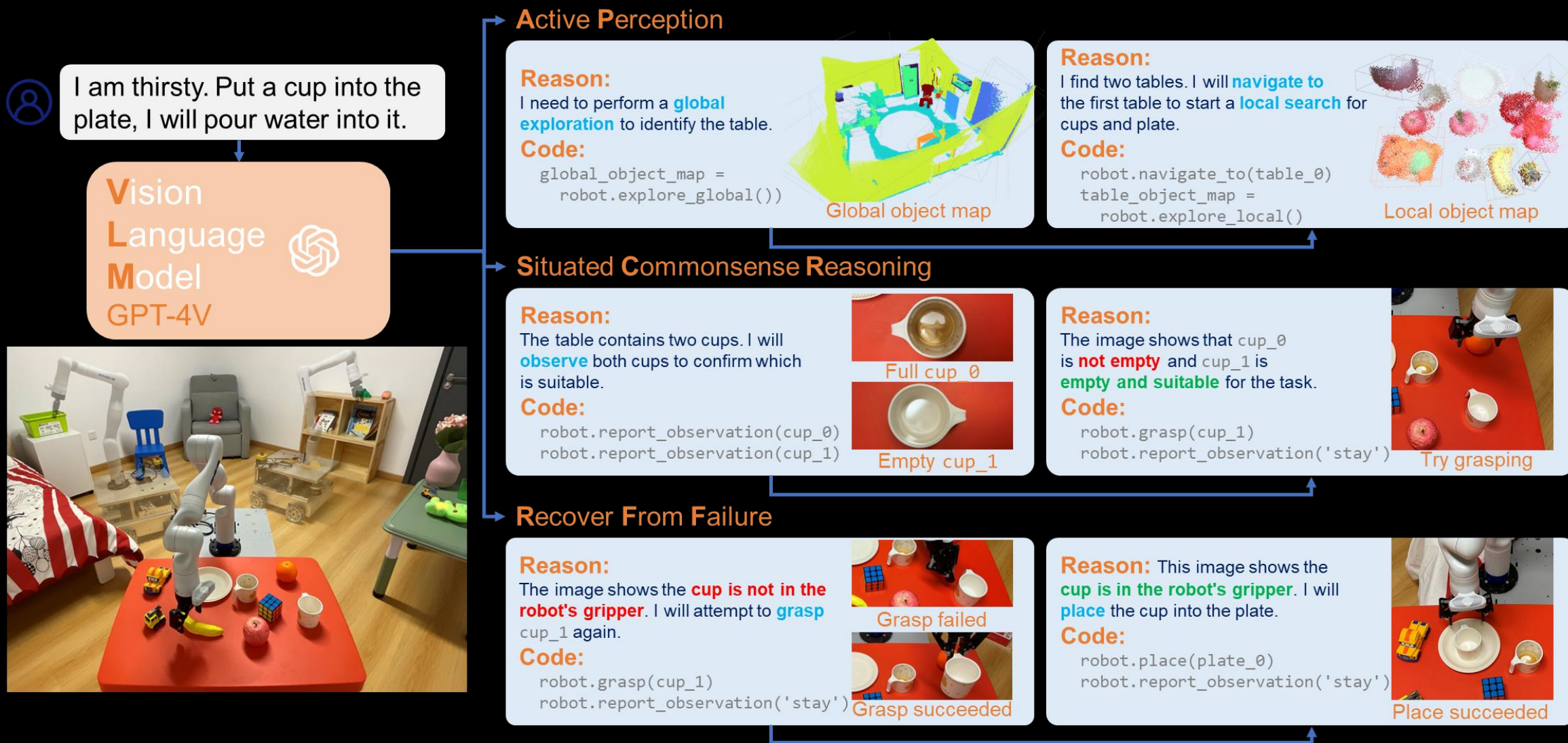
- *(ICRA'25) Closed-Loop Open-Vocabulary Mobile Manipulation with GPT-4V*



Come-Robot



# LLM-based Closed-Loop Open-Vocabulary Manipulation







I am hungry. Could you give me some food? And pass me a cup of juice.

15x

Active  
Perception





**Problem?**

**No learning, just inference**



# Improving Spatial Understanding for VLAs

- (ICML'24) *LEO: An Embodied Generalist Agent in 3D World*
- (ArXiv'25) *LEO-VL: Efficient Scene Representation for Scalable 3D Vision-Language Learning*



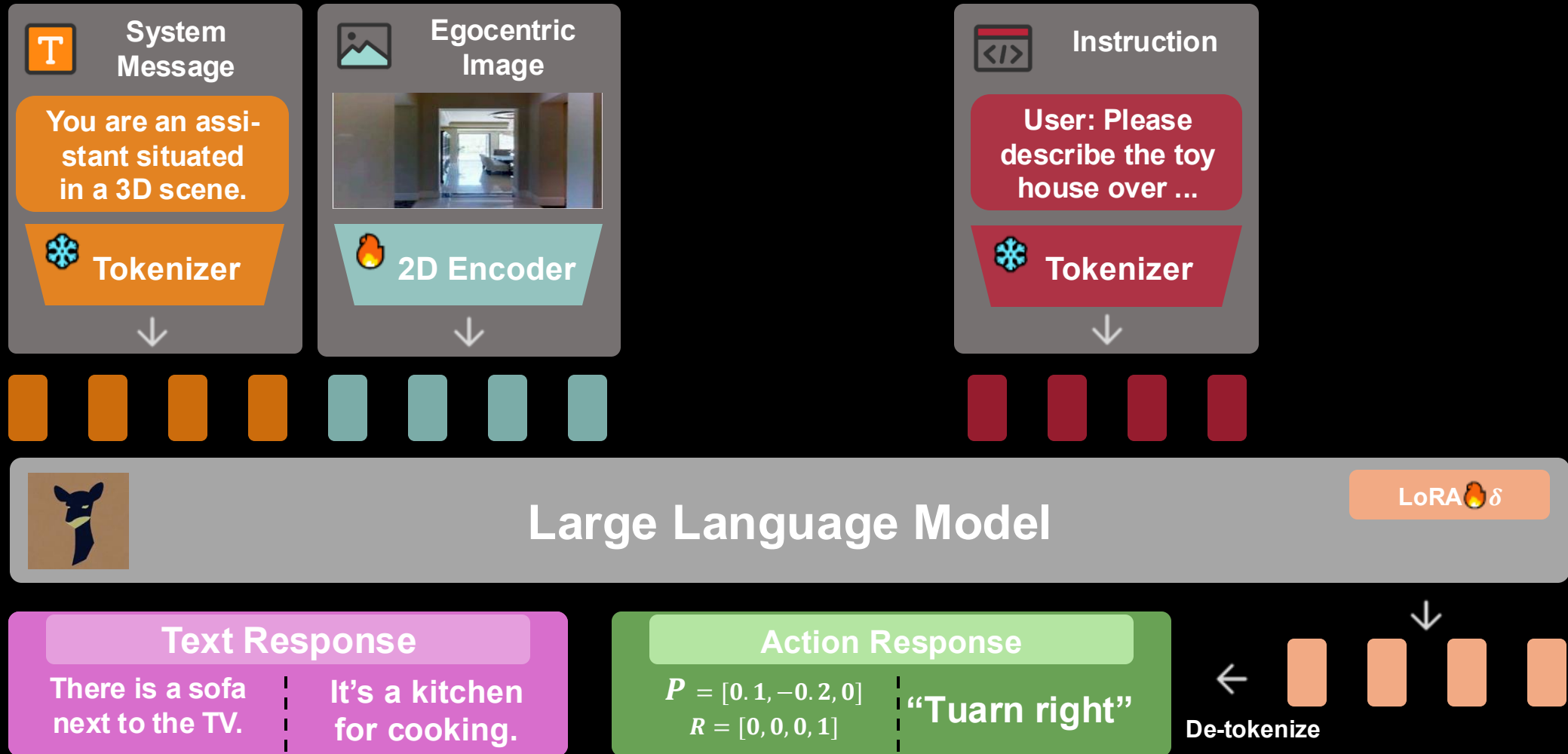
LEO



LEO-VL



# Vision-Language-Action Models





# Can 2D MLLMs Understand 3D Scenes?

## VSI-Bench (CVPR 2025)



### Object Count

How many chairs are there in this room?

Answer: 4

### Relative Distance

Measuring from the closest point of each object, which of these objects (refrigerator, sofa, ceiling light, cutting board) is the closest to the printer?

A. refrigerator B. sofa C. ceiling light D. cutting board

### Appearance Order

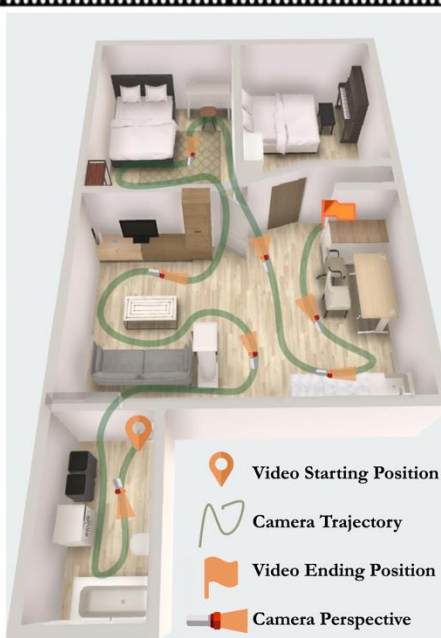
What will be the first-time appearance order of the following categories in the video: basket, printer, refrigerator, kettle?

A. kettle, basket, printer, refrigerator  
B. refrigerator, printer, basket, kettle  
C. basket, printer, refrigerator, kettle  
D. basket, refrigerator, kettle, printer

### Relative Direction

If I am standing by the refrigerator and facing the sofa, is the kettle to my left, right, or back?

A. Left B. right C. back



### Object Size

What is the length of the longest dimension (length, width, or height) of the refrigerator in centimeters?

Answer: 119

### Absolute Distance

Measuring from the closest point of each object, what is the distance between the bed and the sofa in meters?

Answer: 3.2

### Room Size

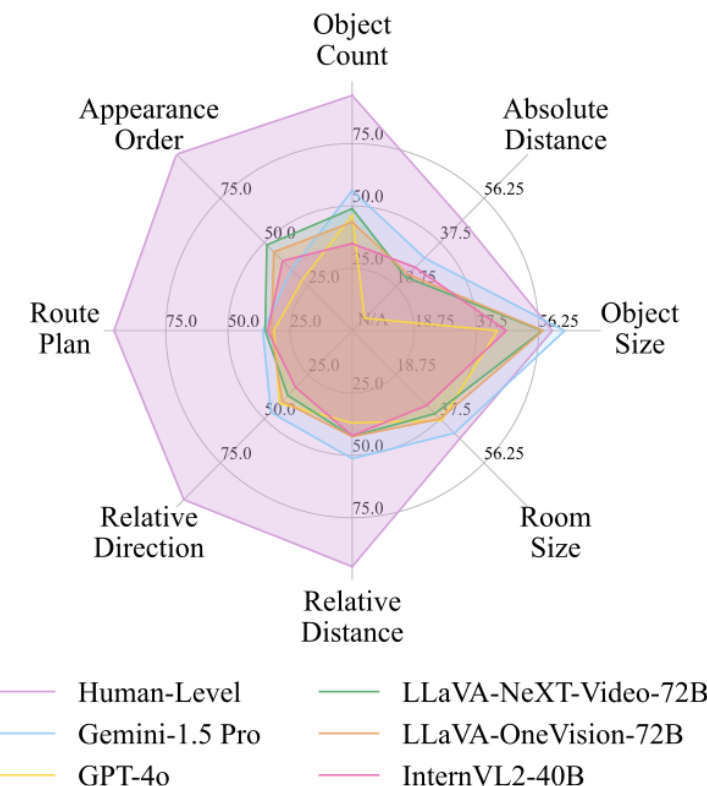
What is the size of this room (in square meters)? If multiple rooms are shown, estimate the size of the combined space.

Answer: 57.6

### Route Plan

You are a robot beginning at the toilet and facing the washer. Navigate to the pan. Fill in this route: 1. Go forward until the washing machine 2. [?] 3. Go forward until the sofa 4. [?] 5. Go forward until the pan.

A. Turn Left, Turn Left B. Turn Left, Turn Right  
C. Turn Back, Turn Right D. Turn Right, Turn Right



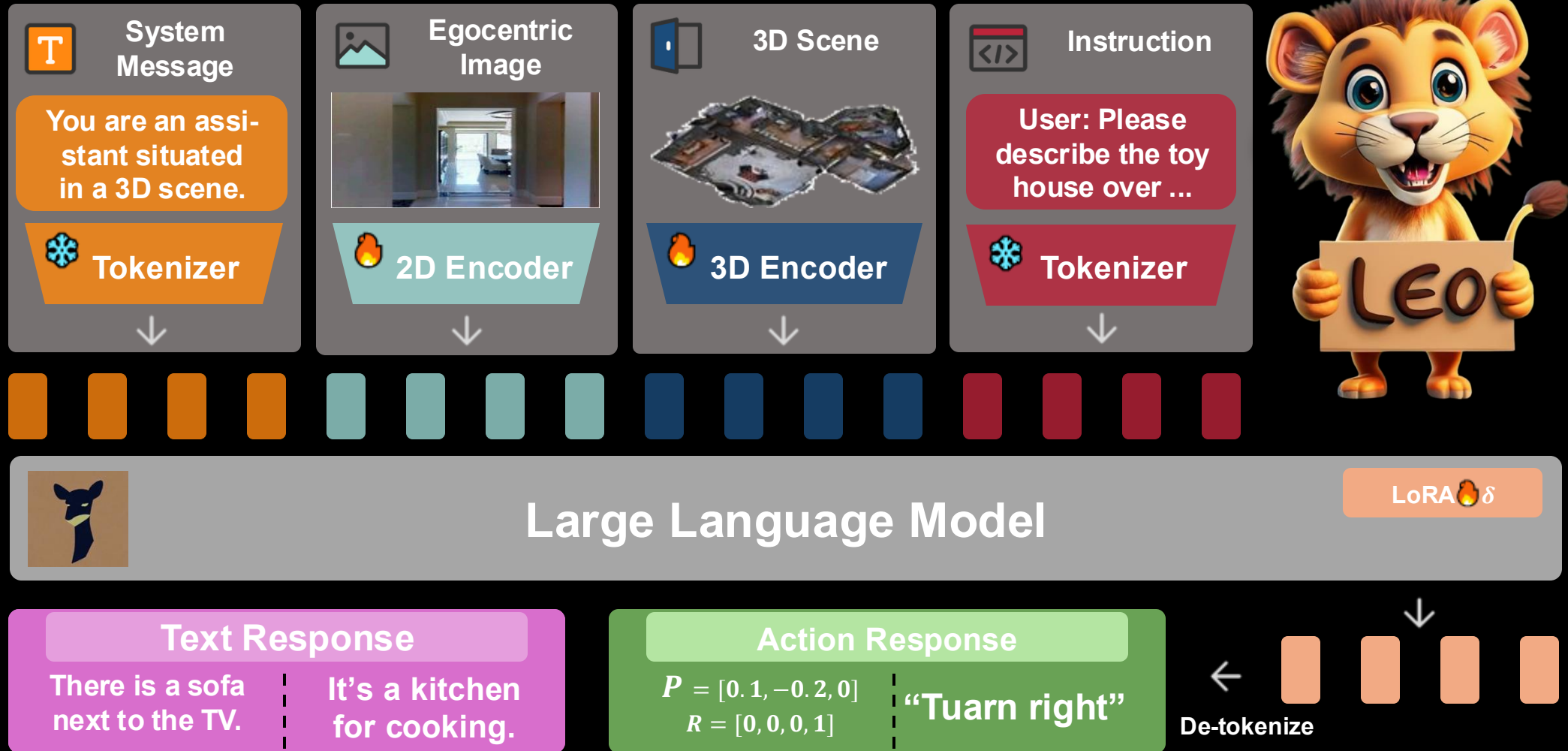
Despite powerful understanding of images and videos, current 2D MLLMs significantly lag far behind humans in **3D scene understanding**, especially **spatial reasoning**.

# Real Deployment?



**Depth is a problem**, but normally not equipped and used by VLAs

# Vision-Language-Action Models





# An Embodied Generalist Agent in 3D World

Jiangyong Huang<sup>1,2\*</sup>, Silong Yong<sup>1,3\*</sup>, Xiaojian Ma<sup>1\*</sup>, Xiongkun Linghu<sup>1\*</sup>, Puhao Li<sup>1,4</sup>,  
Yan Wang<sup>1</sup>, Qing Li<sup>1</sup>, Song-Chun Zhu<sup>1,2,4</sup>, Baoxiong Jia<sup>1</sup>, Siyuan Huang<sup>1</sup>

<sup>1</sup>Beijing Institute for General Artificial Intelligence (BIGAI)

<sup>2</sup>Peking University <sup>3</sup>Carnegie Mellon University <sup>4</sup>Tsinghua University



<https://embodied-generalist.github.io/>

## Embodied Generalist Agent

Capabilities: *Perception*, *Grounding*, *Reasoning*, *Planning*, *Acting*

### Tasks

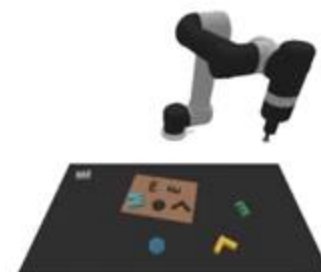
3D Object Captioning  
Scene Captioning

3D Question Answering  
Embodied Reasoning

3D Dialogue  
Task Planning

Embodied Navigation  
Robotic Manipulation

### 3D World



# Efficient Representation Bridging 2D-3D Perception

Pros

Cons

3D perception

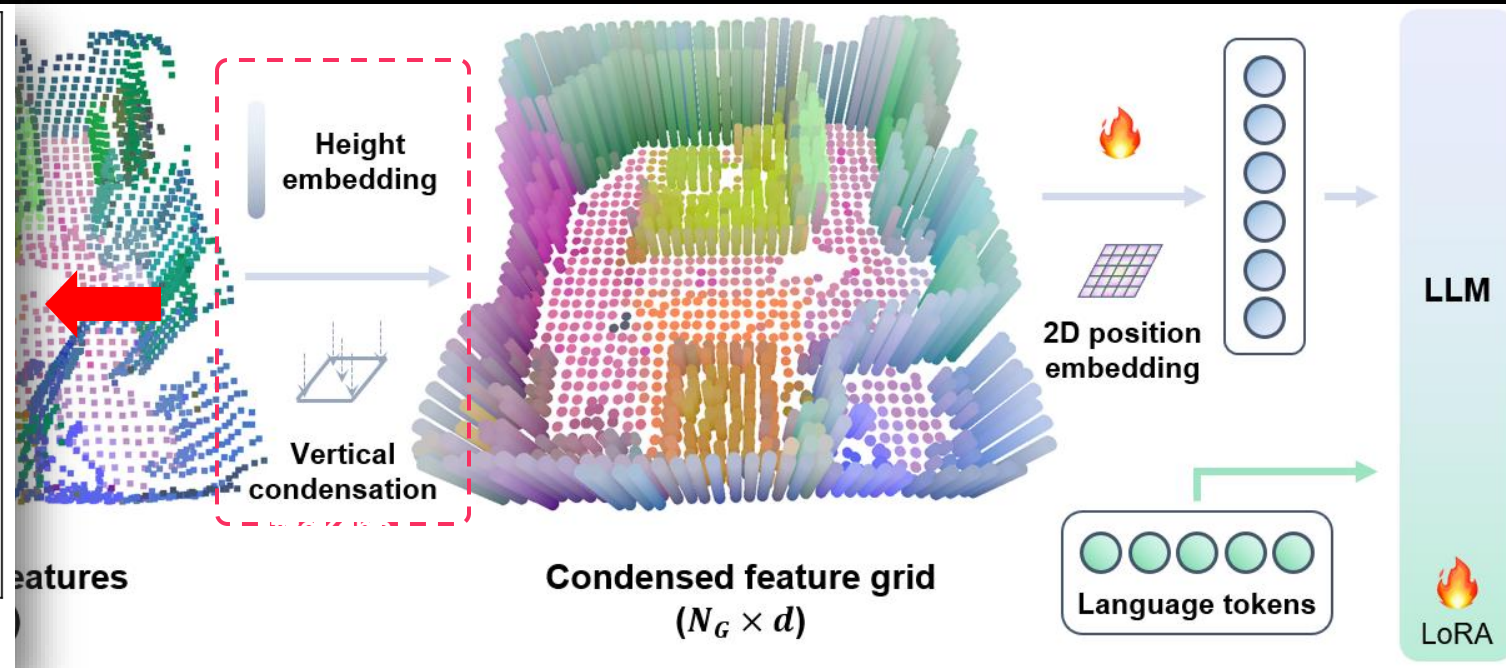
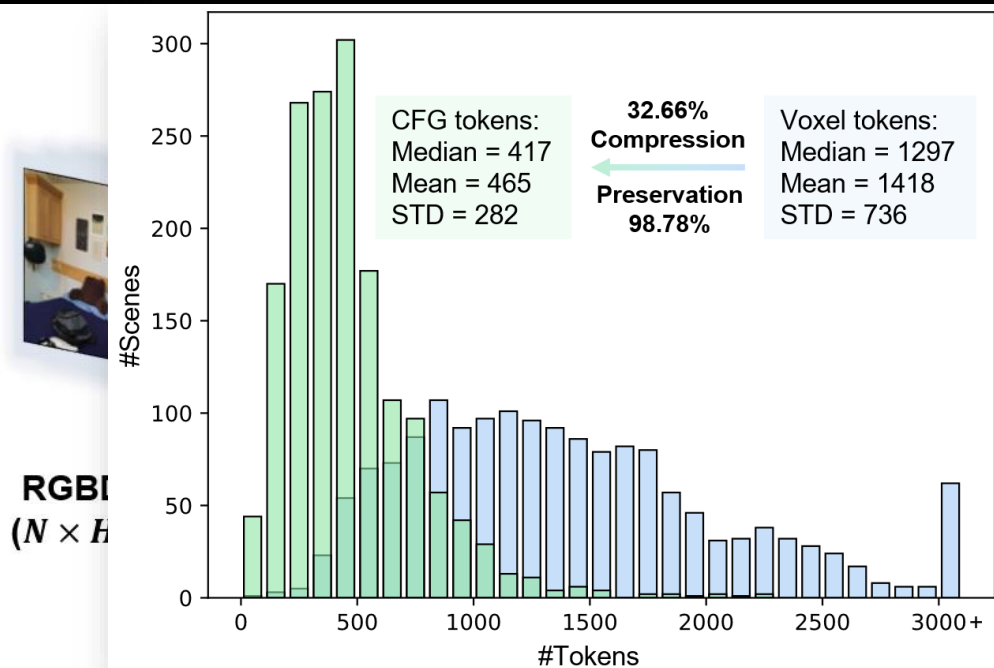
Explicit 3D structure

Complex pre-processing pipelines, learning difficulty

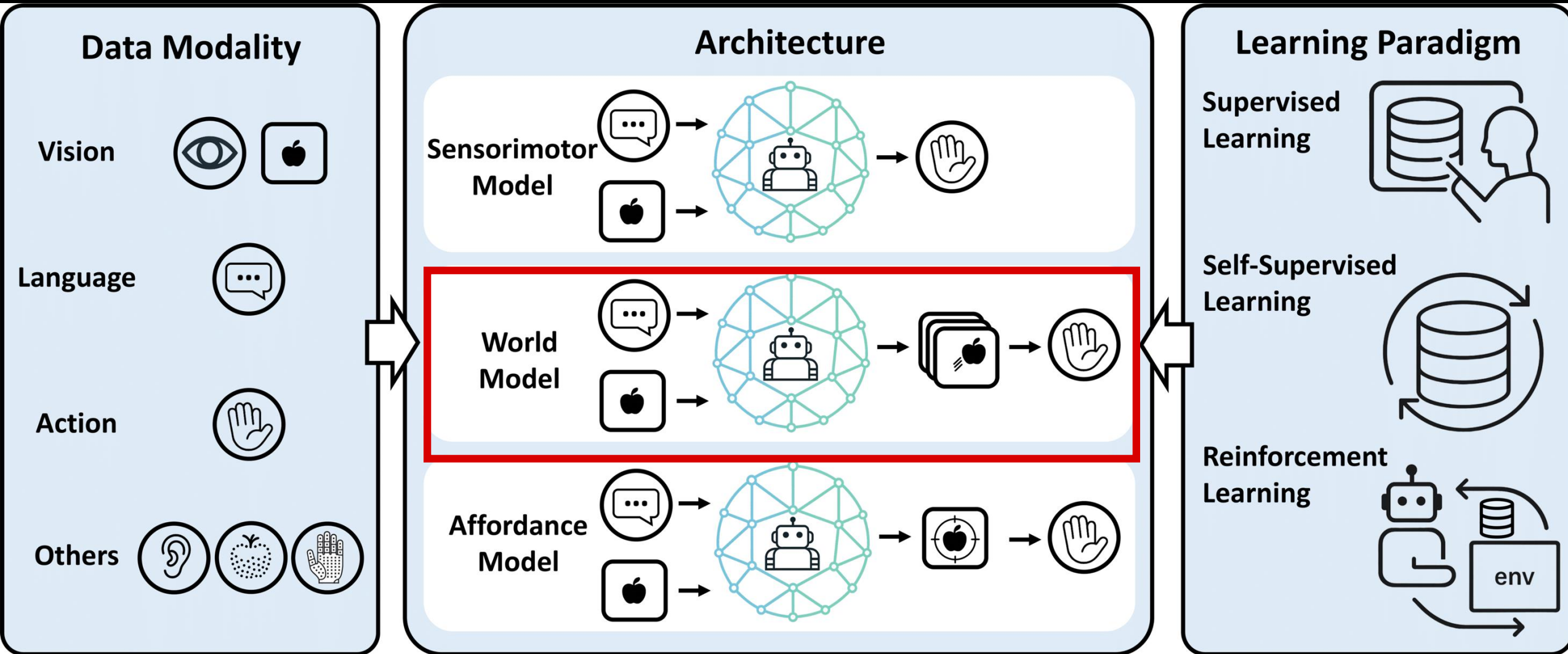
2D perception

Strong capability

Significant computation overhead (thousands of tokens)



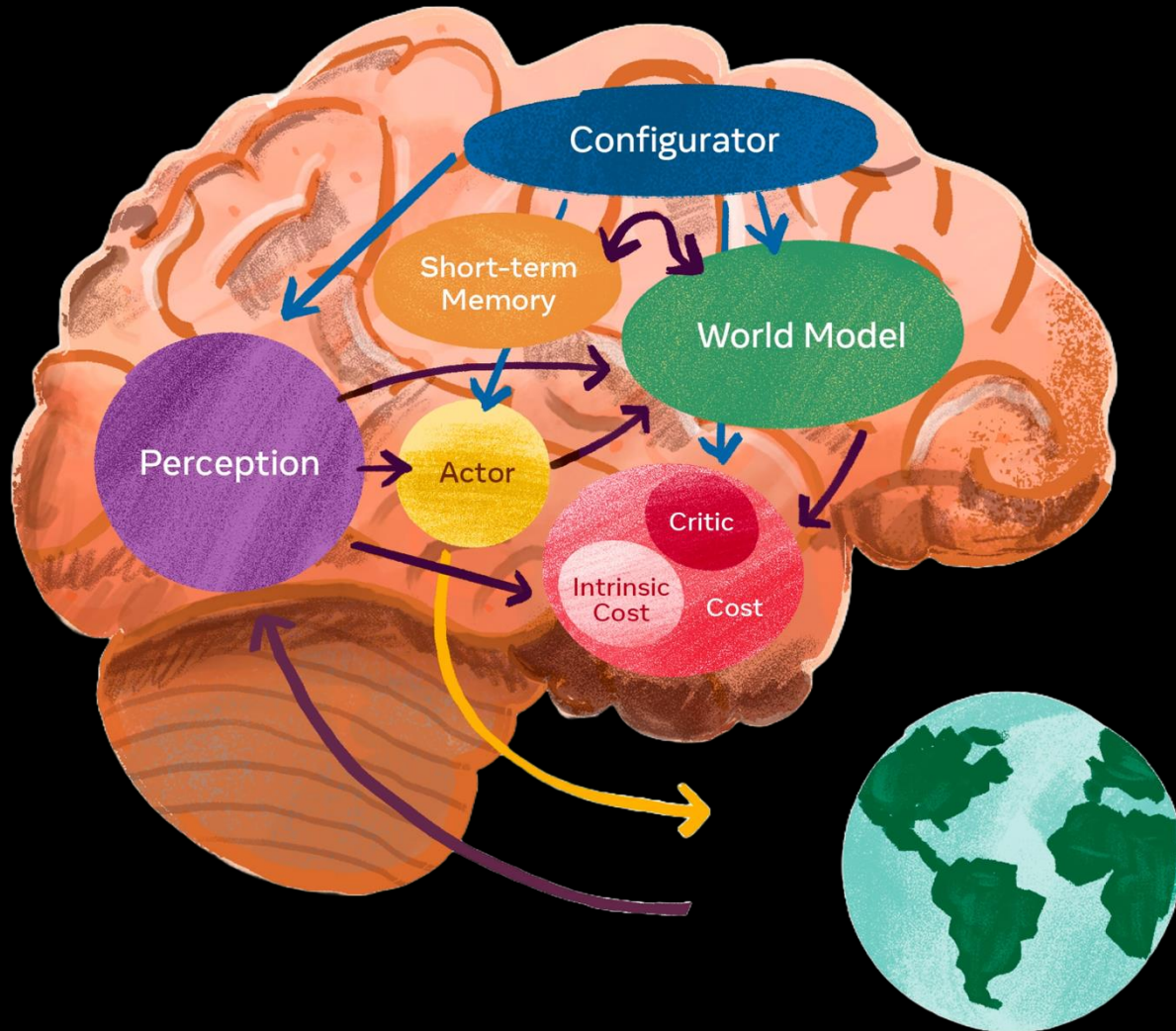
# A Closer Look at VLA Model Design



Vision-Language-Action Models for Robotics (IEEE Access 2025)



# World Models



*“If the organism carries a **small-scale model** of external reality and of its own possible actions within its head, it is able to **try out various alternatives**, conclude which is the best of them, **react to future situations before they arise**, utilize the knowledge of past events in dealing with the present and future.”*

— Kenneth Craik (1943)

0 min

## Model-based RL

Representation learning for long-horizon tasks

Under game setting

Dreamer 4, Google DeepMind 2025

## Video Generation

Flexible conditional generation

Weak physical consistency / modeling of action

Veo 3.1, Google Deepmind 2025

## Latent Action Learning

Aligning video generation with latent actions

Limited by the view-point

DreamGen, NVIDIA GEAR 2025

## Spatial Representations

World modeling with 3D Gaussians

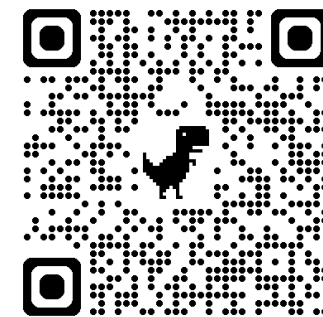
Interactiveness for robot manipulation?

Marble, WorldLabs 2025

# Scalable World Modeling with 3D Gaussians

- *(ICCV'25) GWM: Towards Scalable Gaussian World Modeling for Robotic Manipulation*

GWM





# Encoding 3D Gaussians into Latent Space



(Optional)

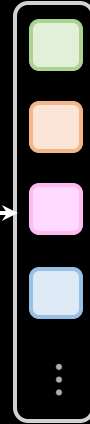
Unposed Img.

Splatt3R



Gaussian Splats  $\mathcal{G}_t$

3D VAE



Compact Latent Representation

Pos.  
Emb.

3D VAE



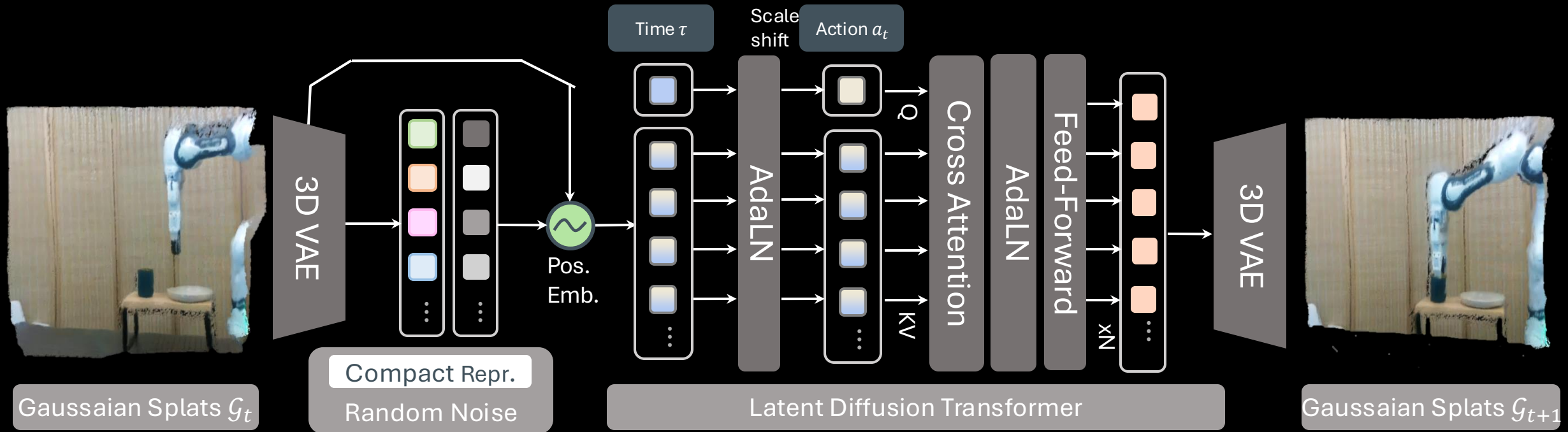
Gaussian Splats  $\mathcal{G}_t$

**Feed-Forward 3D  
Gaussian Reconstruction**

**FPS-based Subsampling  
Query-based Encoding**

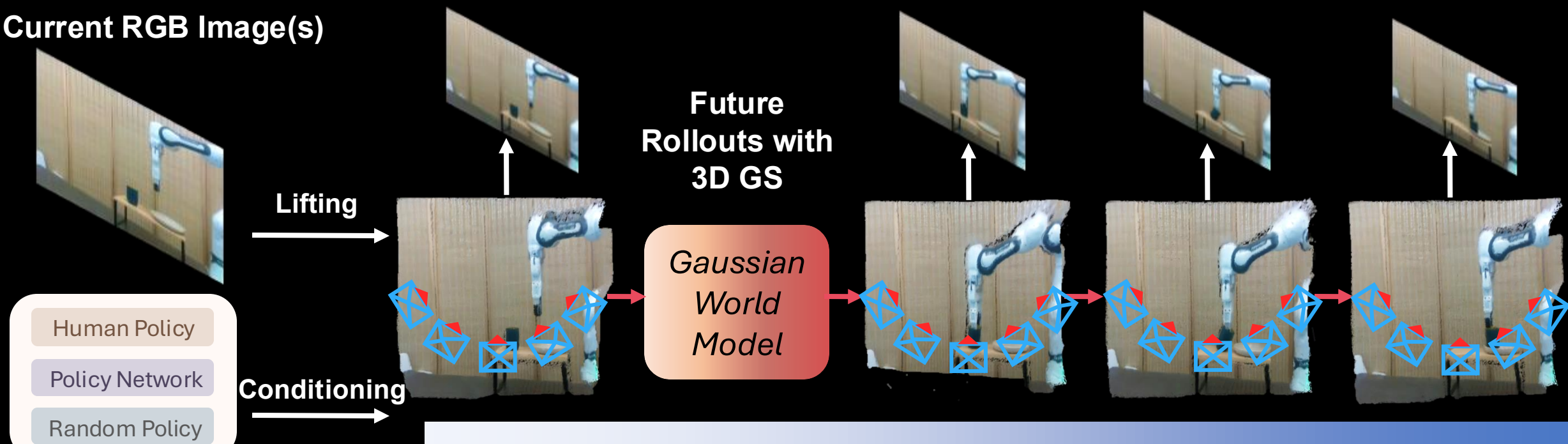
**Rendering / Geometry  
Supervision**

# GWM: Gaussian World Model

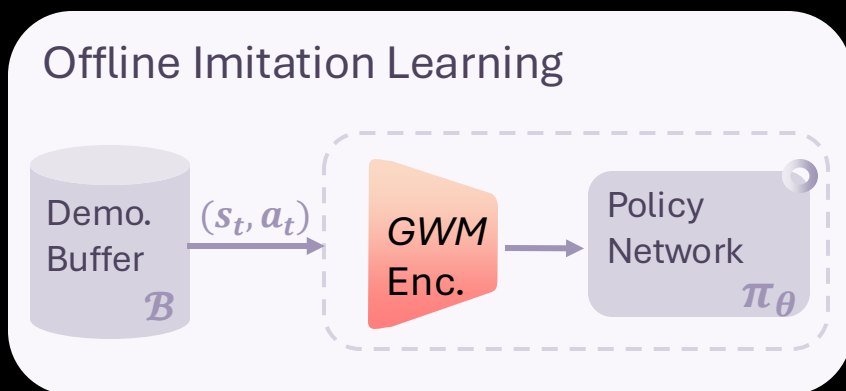


**DiT-based Dynamics Learning and Prediction**

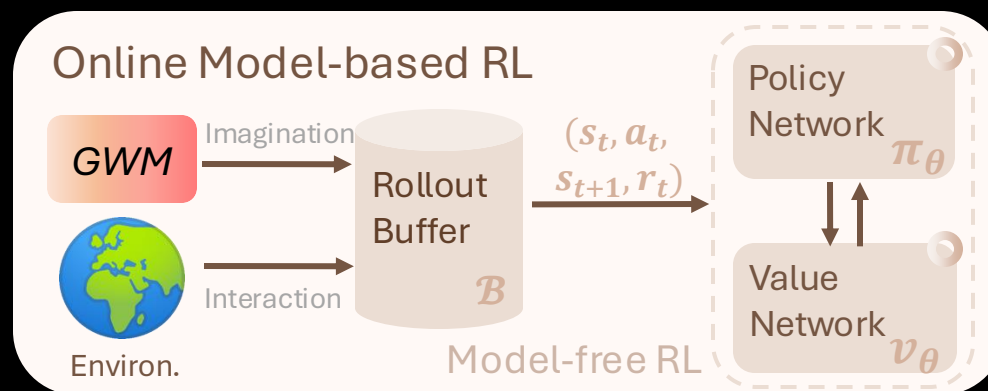
## Current RGB Image(s)



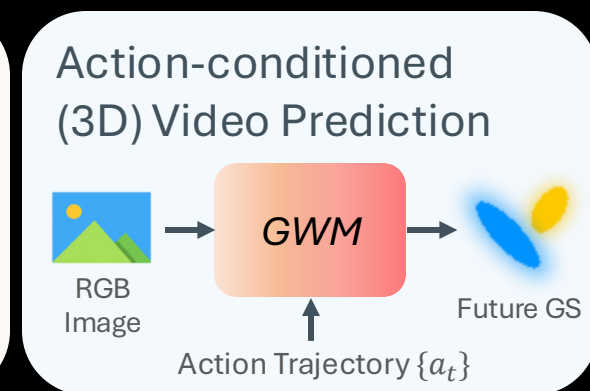
## Offline Imitation Learning



## Online Model-based RL

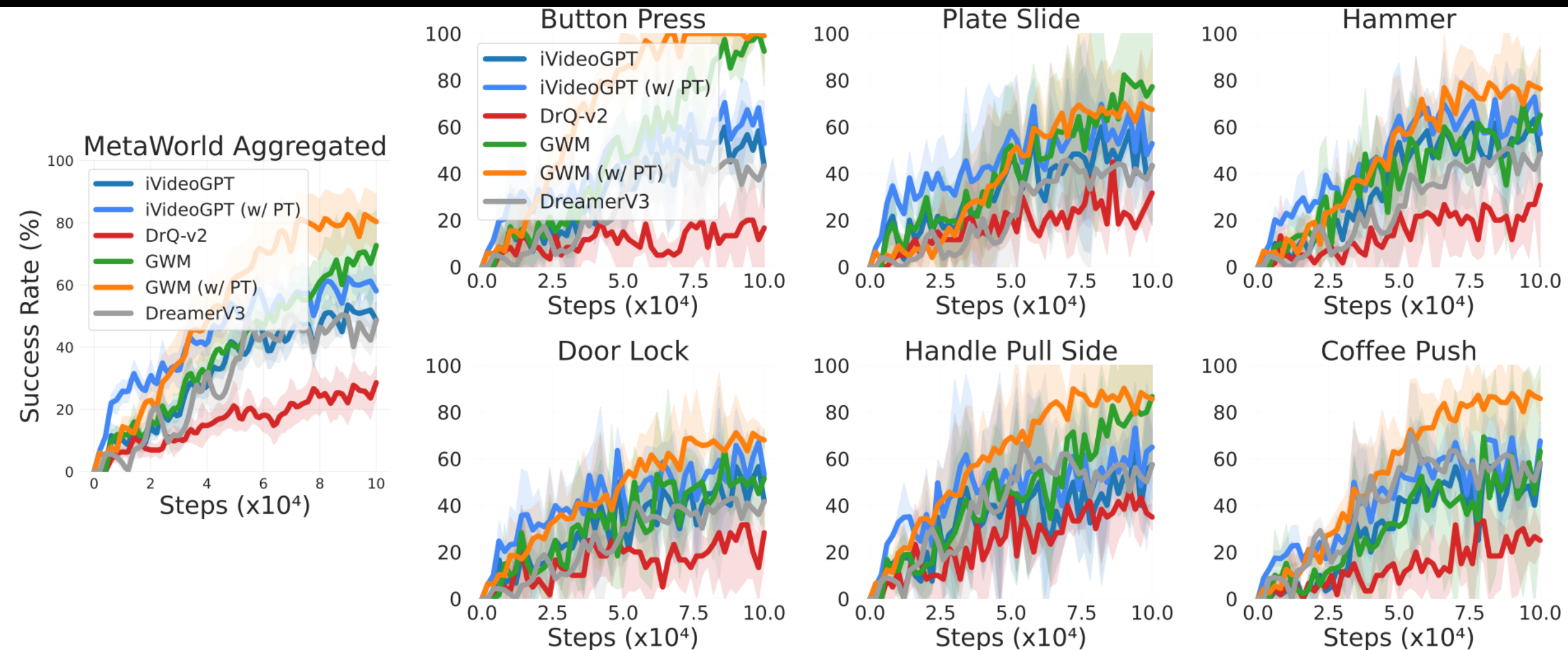


## Action-conditioned (3D) Video Prediction





# GWM for Online Model-based RL



**Additional reward learning on top of GWM for online RL**

# GWM for Real-World Robot Manipulation



Diffusion  
Policy

DP w/  
GWM

Comparison



| FRANKA-PNP       | Diffusion Policy | GWM (Ours)   |
|------------------|------------------|--------------|
| Cup distractor   | 6/10             | 7/10         |
| Plate distractor | 1/5              | 3/5          |
| Table distractor | 0/5              | 3/5          |
| <b>Total</b>     | <b>7/20</b>      | <b>13/20</b> |



### 3. Evolving Force-Aware Control Skills for Human-Robot Interaction



# Can Humanoids Interact at This Level?

Humans effortlessly **squat to retrieve objects** from the ground and then **walk to another distant place**.



Enable **holistic** and **long-horizon** humanoid–scene interaction

# Whole-body demonstrations is significantly limited

GPT/Qwen

1.2B Hours

$\pi_0$

10k Hours

autonomous 1

Gr00T

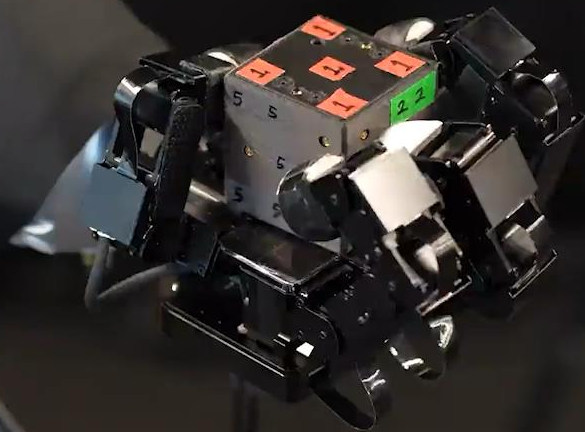
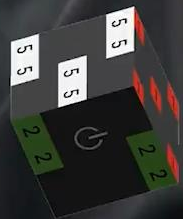
88 Hours

Whole-body data

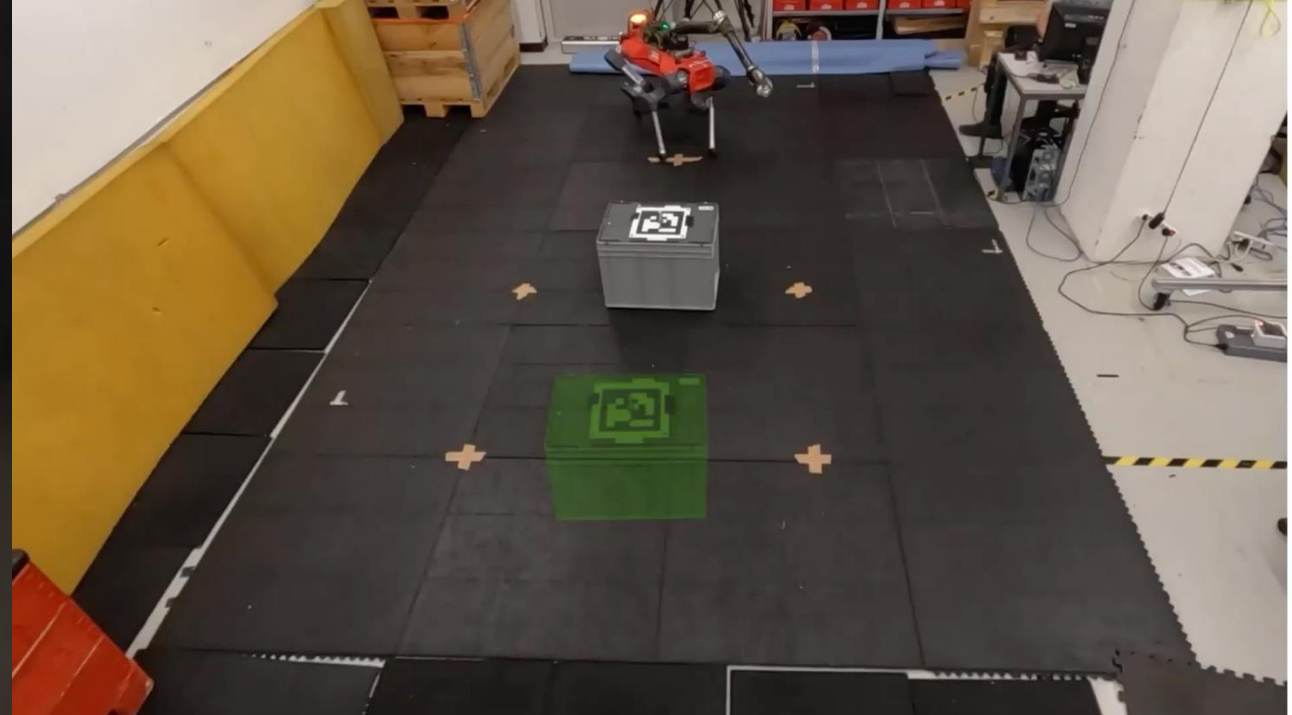
~0 hours



Goal



1x



Task 1  
Mustard Bottle  
2x Speed





# Goal of Action and Control

- Enable **agile and stable** whole-body control for humanoid robots
- Mitigate the missing **force modality** for contact-rich manipulation tasks
- Safe and helpful **human-robot collaboration** patterns

# Force-Aware Manipulation

UniFP (CoRL'25)

Best Paper

# Teleoperation / Tracking

CLONE (CoRL'25)

Speed Tracking, 1.25x Speed

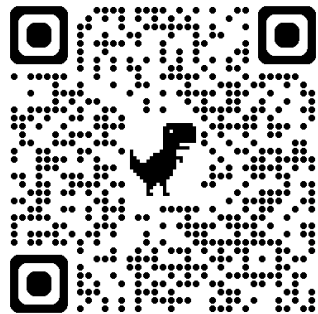
# Human-Robot Interaction

COLA (ArXiv'25)

# Agile Humanoid Whole-Body Teleoperation

- (CoRL'25) *CLONE*: Closed-Loop Whole-Body Humanoid Teleoperation for Long-Horizon Tasks

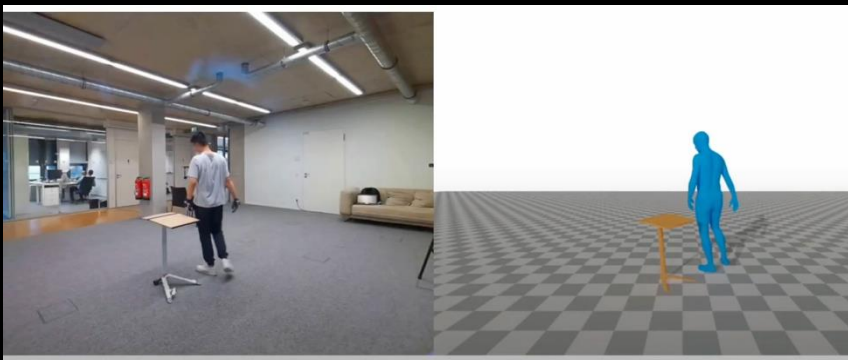
CLONE



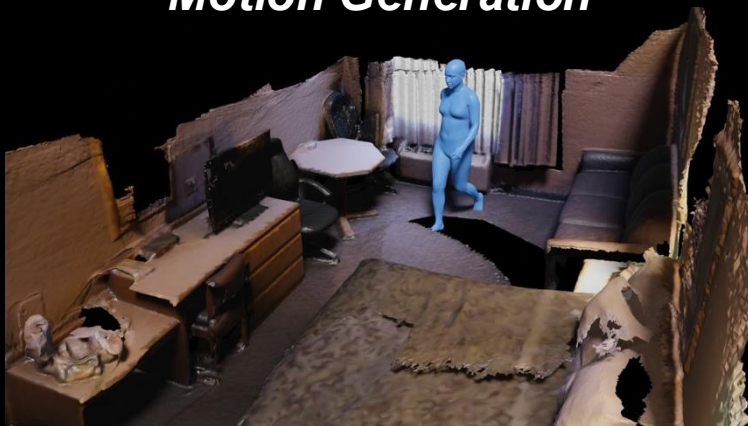


# Driving a Physical Humanoid with Human Motion

*MoCap*



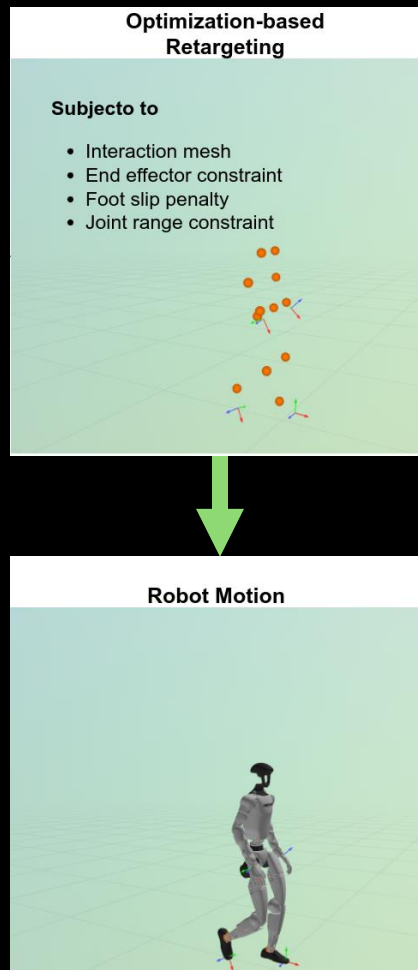
*Motion Generation*



*Human Motion*



*LAFAN  
Motion Retargeting*



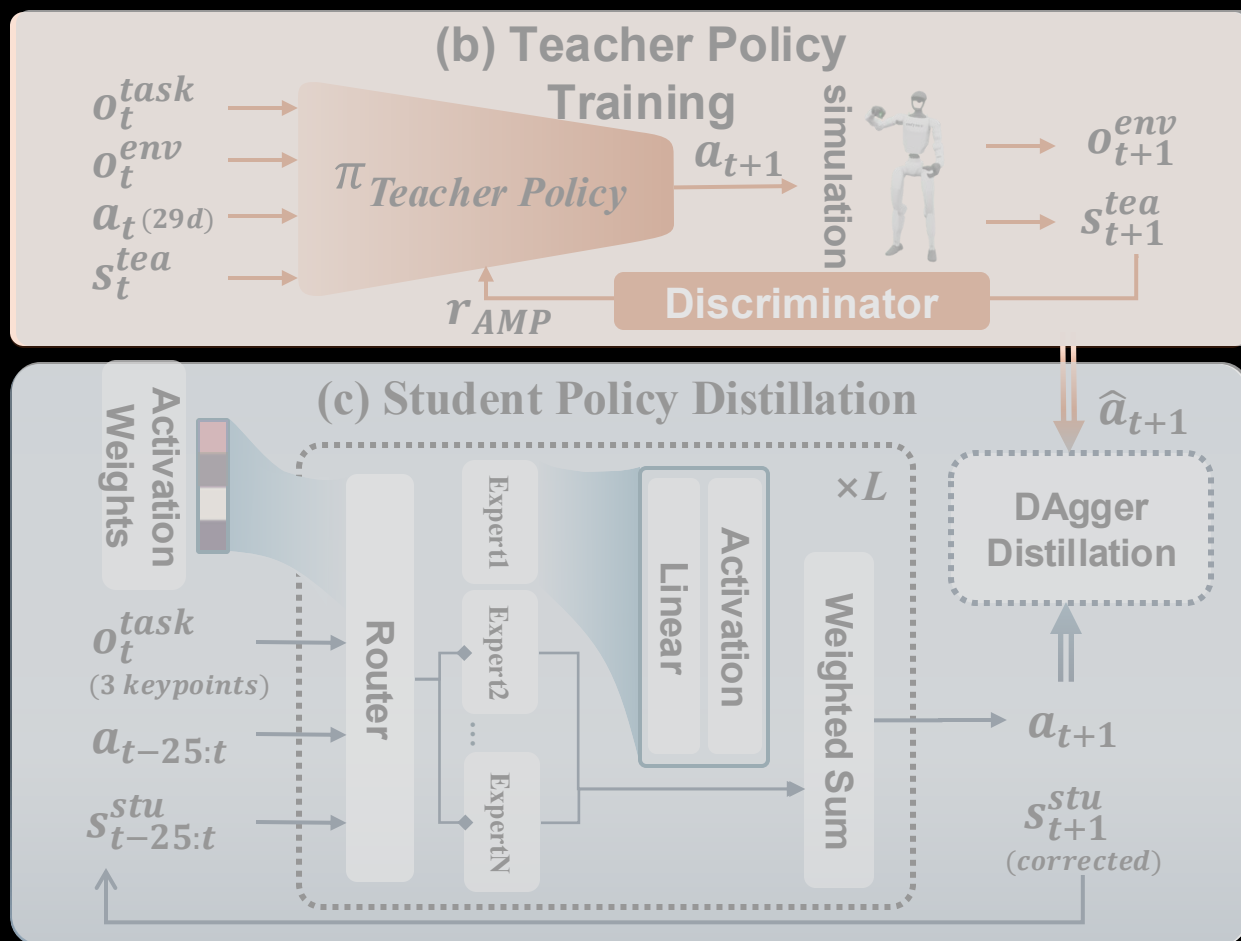
*Humanoid  
Controller*

MLP  
MoE  
Transformer  
PID  
MPC  
...



*Robot Execution*

# CLONE: Humanoid Whole-Body Teleoperation



Learning a teacher policy with privileged information for human motion tracking

Distilling a MoE-based student policy with Behavior Cloning (Dagger)



**BIGAI**



**UNITREE**

# CLONE

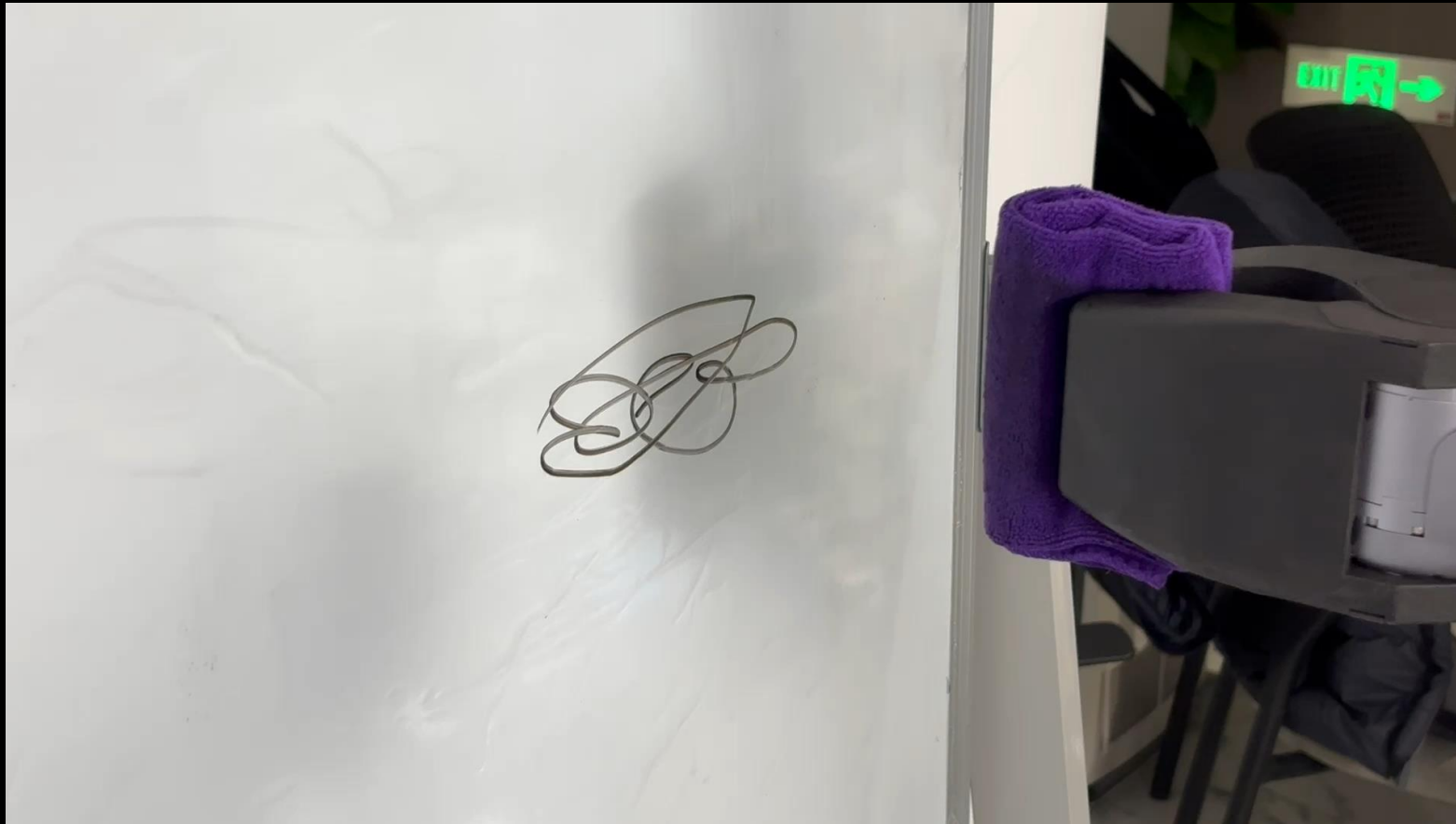
**Holistic Closed-Loop Whole-Body  
Humanoid Teleoperation for Long-Horizon Tasks**





OK, can we let the robot help us wipe the whiteboard first after meeting?

Let me **collect the data** and **imitation learning** will solve the rest 😊

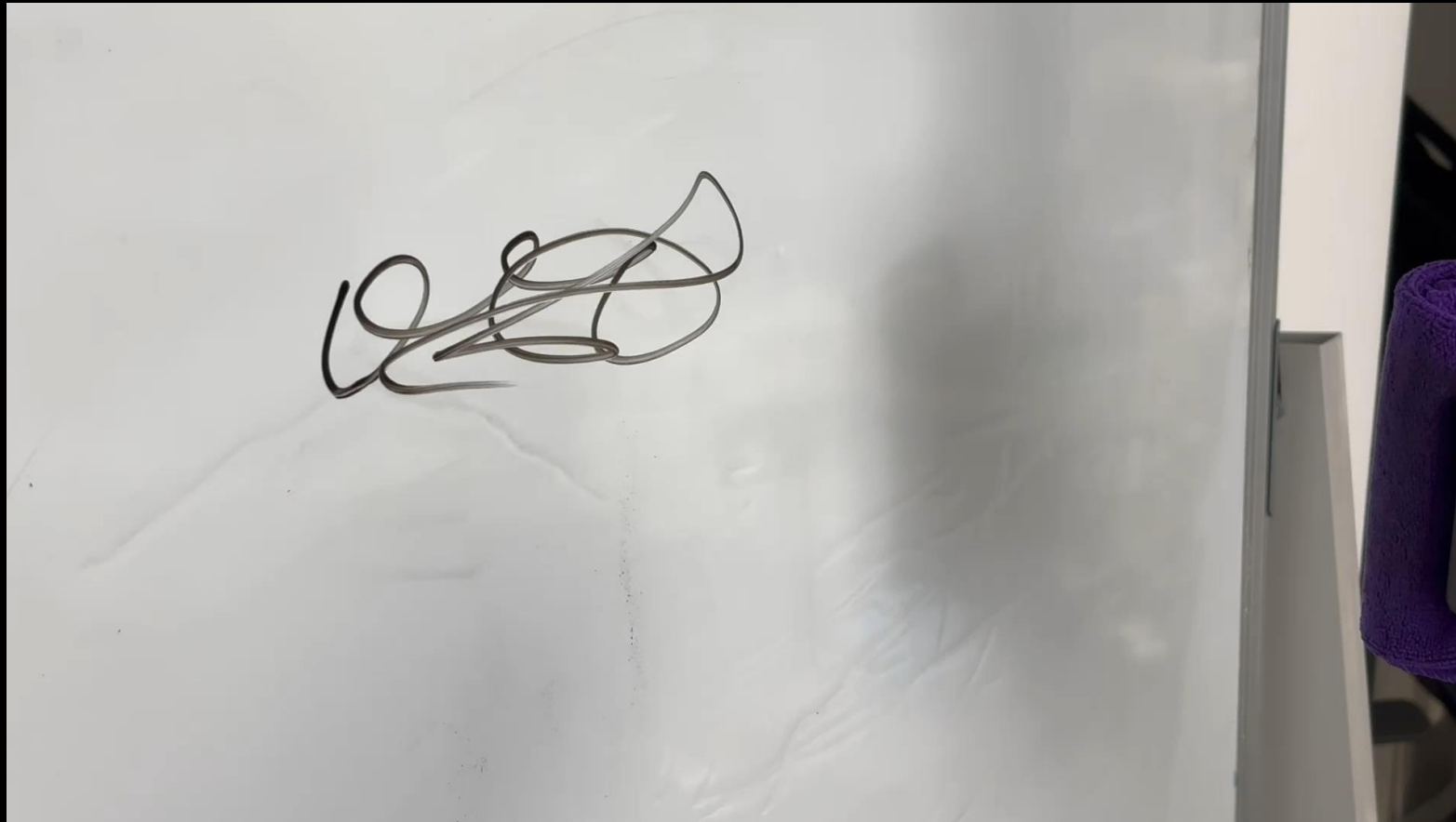


Data collection is a **disaster**

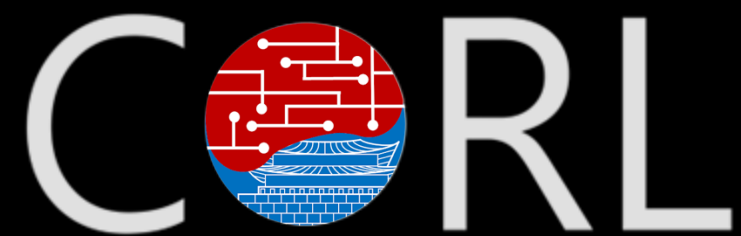


OK, can we let the robot help us wipe the whiteboard first after meeting?

Let me **collect the data** and **imitation learning** will solve the rest 😊



Of course, the learned policy **failed no matter how much data used** 😞

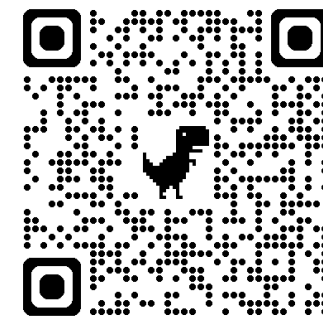


# Unified Force-Position Control Policy

- (CoRL'25 *Best Paper*) *Learning a Unified Policy for Position and Force Control in Legged Loco-Manipulation*



UniFP





# Revisiting the control formulation

mass-spring-damper system

$$\mathbf{F} = K(\mathbf{x} - \boxed{\mathbf{x}^{\text{cmd}}}) + \cancel{D(\dot{\mathbf{x}} - \boxed{\dot{\mathbf{x}}^{\text{cmd}}})} + \cancel{M(\ddot{\mathbf{x}} - \boxed{\ddot{\mathbf{x}}^{\text{cmd}}})}$$

$$\mathbf{x} = \mathbf{x}^{\text{cmd}} + \frac{\mathbf{F}}{K}$$

And if the end effector moves really slowly...

# Revisiting the control formulation

mass-spring-damper system

$$\mathbf{F} = K(\mathbf{x} - \mathbf{x}^{\text{cmd}}) + \cancel{D(\dot{\mathbf{x}} - \dot{\mathbf{x}}^{\text{cmd}})} + \cancel{M(\ddot{\mathbf{x}} - \ddot{\mathbf{x}}^{\text{cmd}})}$$

$$\mathbf{x} = \mathbf{x}^{\text{cmd}} + \frac{\mathbf{F}}{K}$$

Force can be estimated via **position offsets!**

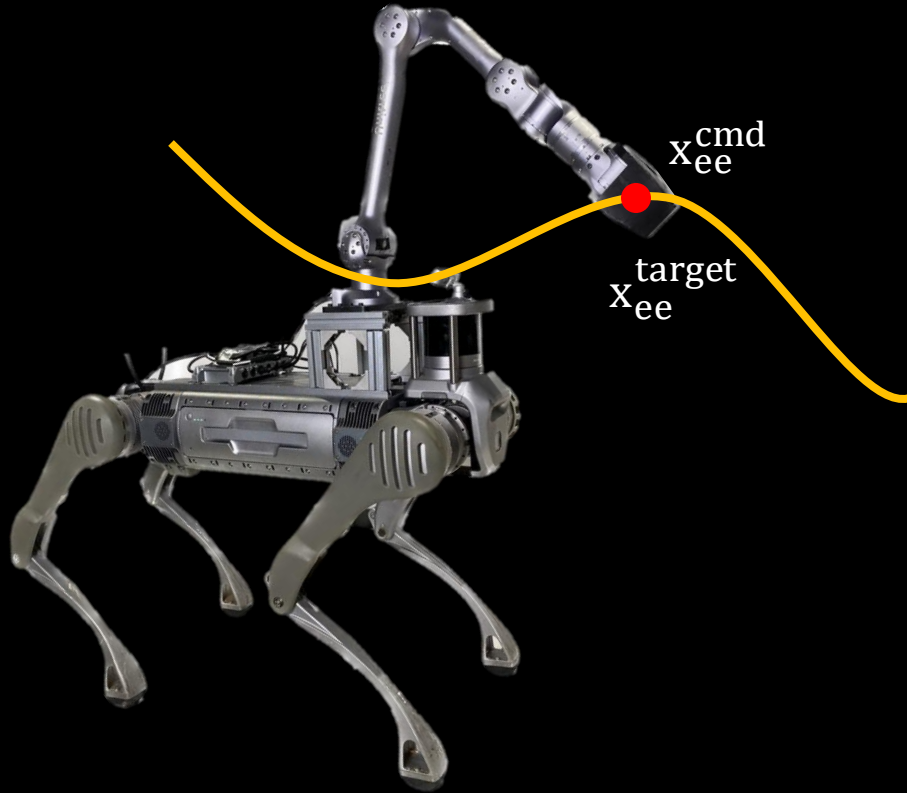
Tracking the **force-adjusted position** enables joint force-position control.

# Formulating forces with positions

$$F = K(x - x^{\text{cmd}})$$

## Position control

$$x^{\text{target}} = x^{\text{cmd}}$$

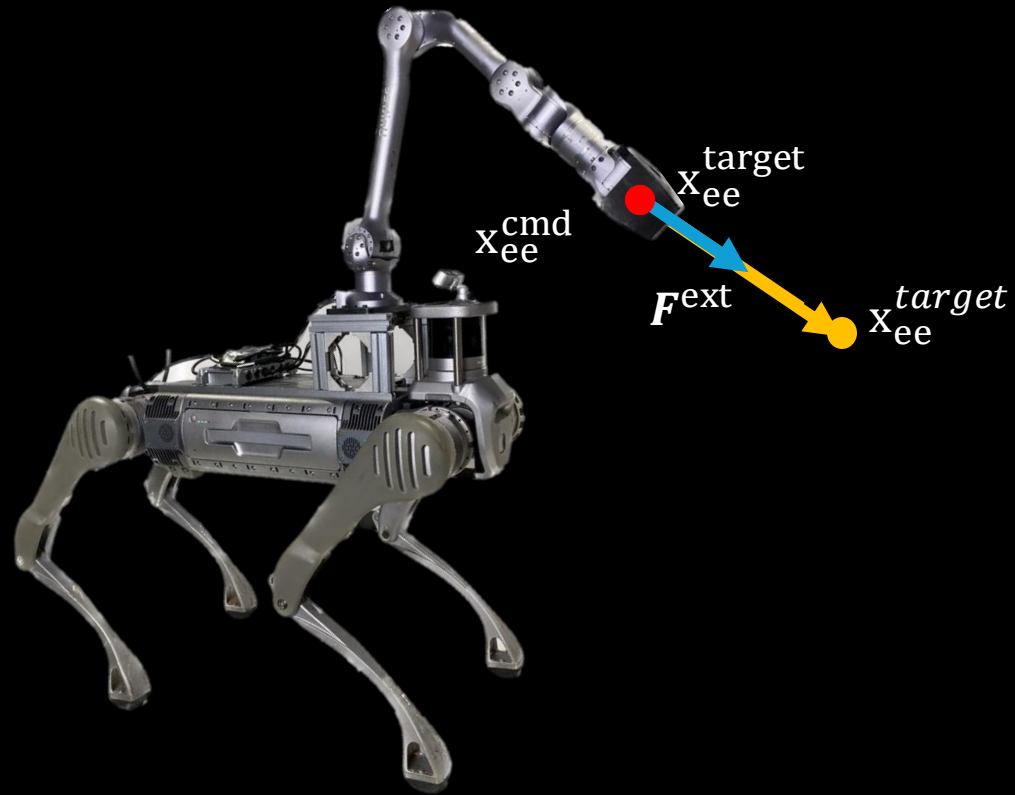




# Formulating forces with positions

$$F = K(x - x^{\text{cmd}})$$

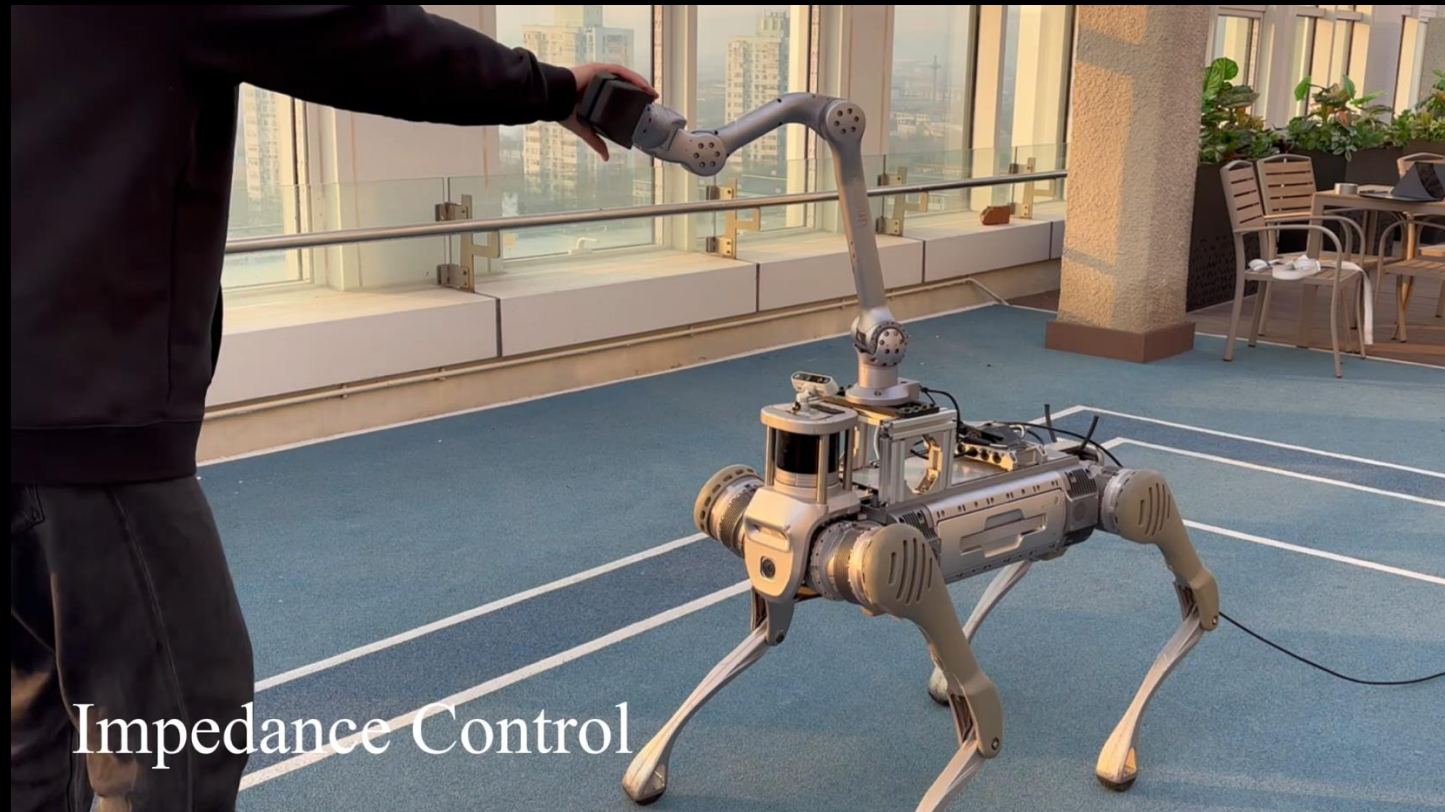
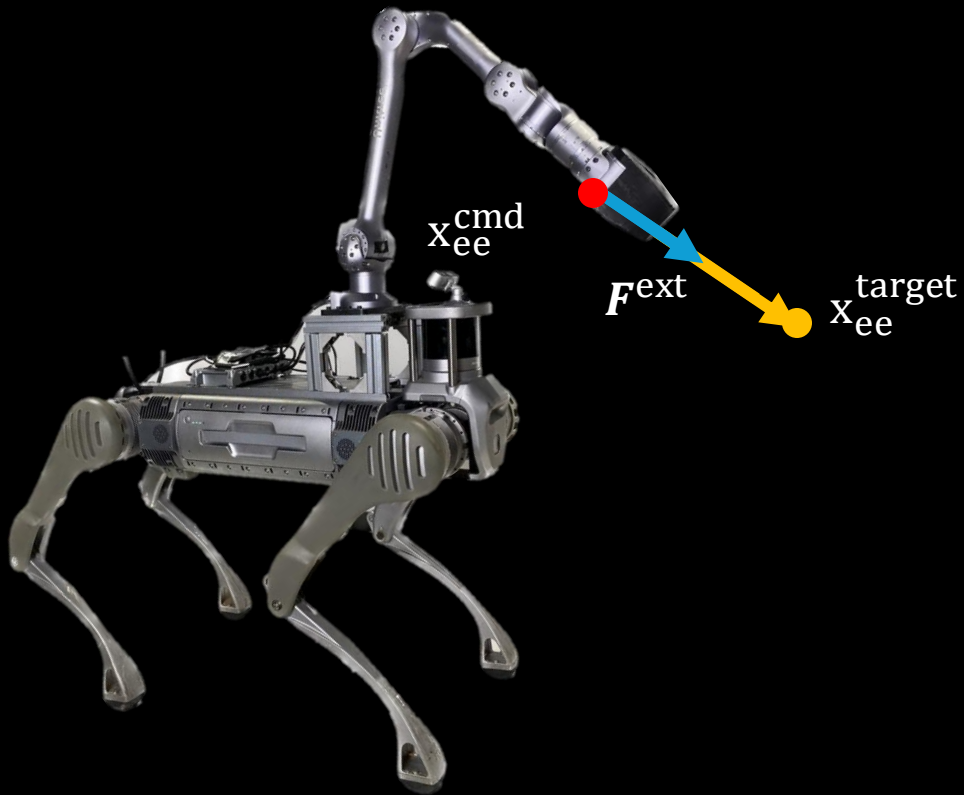
When with external force



# Formulating forces with positions

$$\mathbf{F} = K(\mathbf{x} - \mathbf{x}^{\text{cmd}})$$

Impedance control  $\mathbf{x}^{\text{target}} = \mathbf{x}^{\text{cmd}} + \frac{\mathbf{F}^{\text{ext}}}{K}$



# Revisiting the control formulation

mass-spring-damper system

$$\mathbf{F} = K(\mathbf{x} - \mathbf{x}^{\text{cmd}}) + D(\dot{\mathbf{x}} - \dot{\mathbf{x}}^{\text{cmd}}) + M(\ddot{\mathbf{x}} - \ddot{\mathbf{x}}^{\text{cmd}})$$

And if we care about the locomotion



# Revisiting the control formulation

mass-spring-damper system

$$\mathbf{F} = K(\mathbf{x} - \mathbf{x}^{\text{cmd}}) + D(\dot{\mathbf{x}} - \dot{\mathbf{x}}^{\text{cmd}}) + M(\ddot{\mathbf{x}} - \ddot{\mathbf{x}}^{\text{cmd}})$$

$$\dot{\mathbf{x}} = \dot{\mathbf{x}}^{\text{cmd}} + \frac{\mathbf{F}}{D}$$

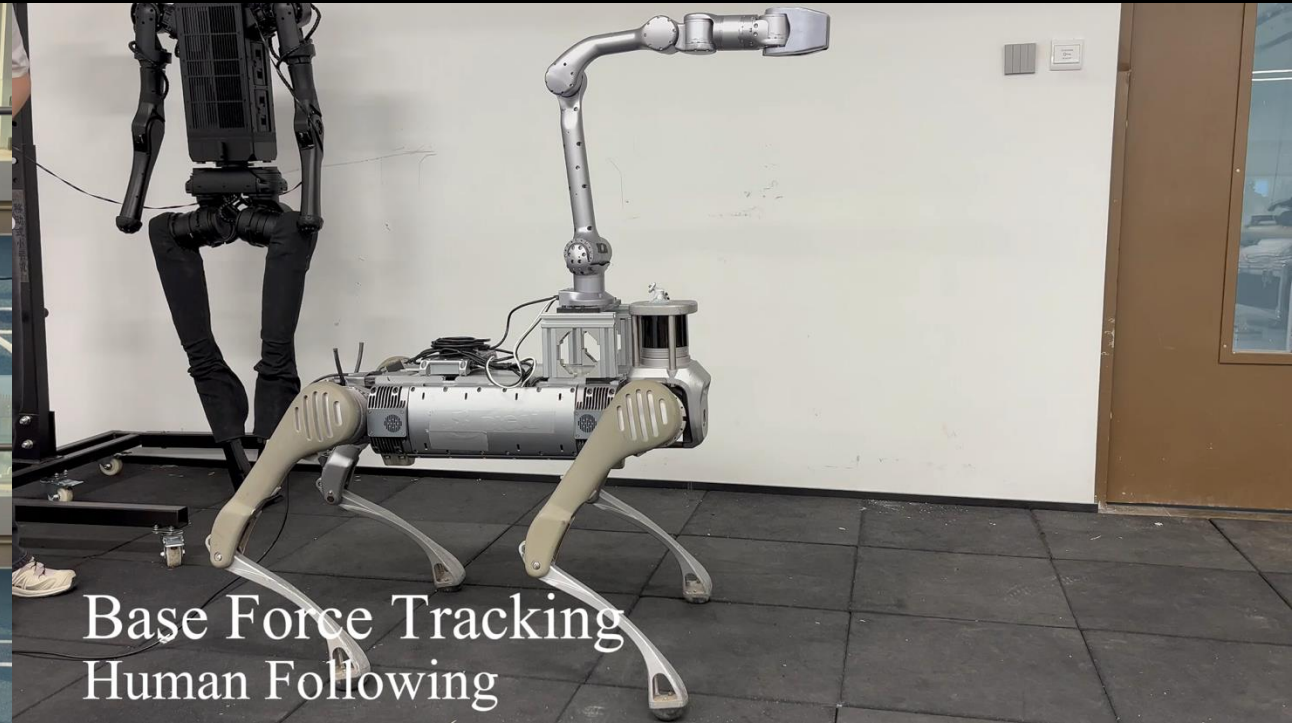
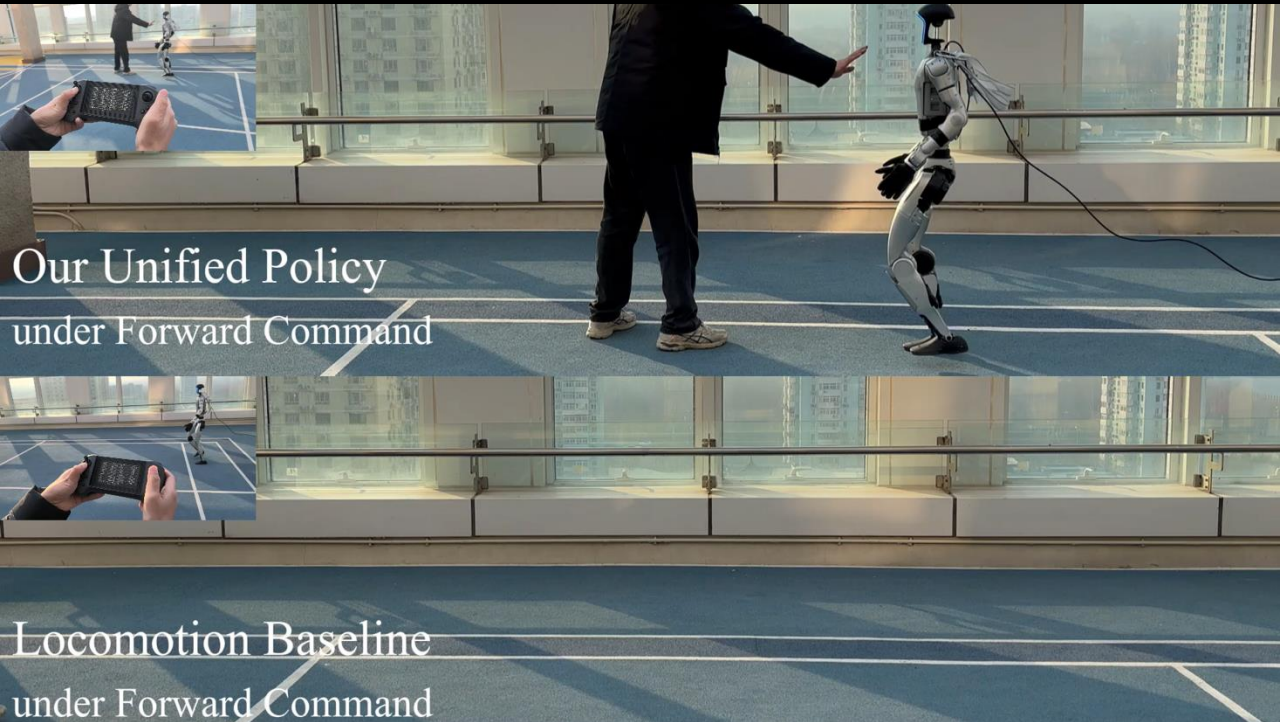
**Force-adjusted velocity** enables compliant locomotion

# Formulating forces with velocities

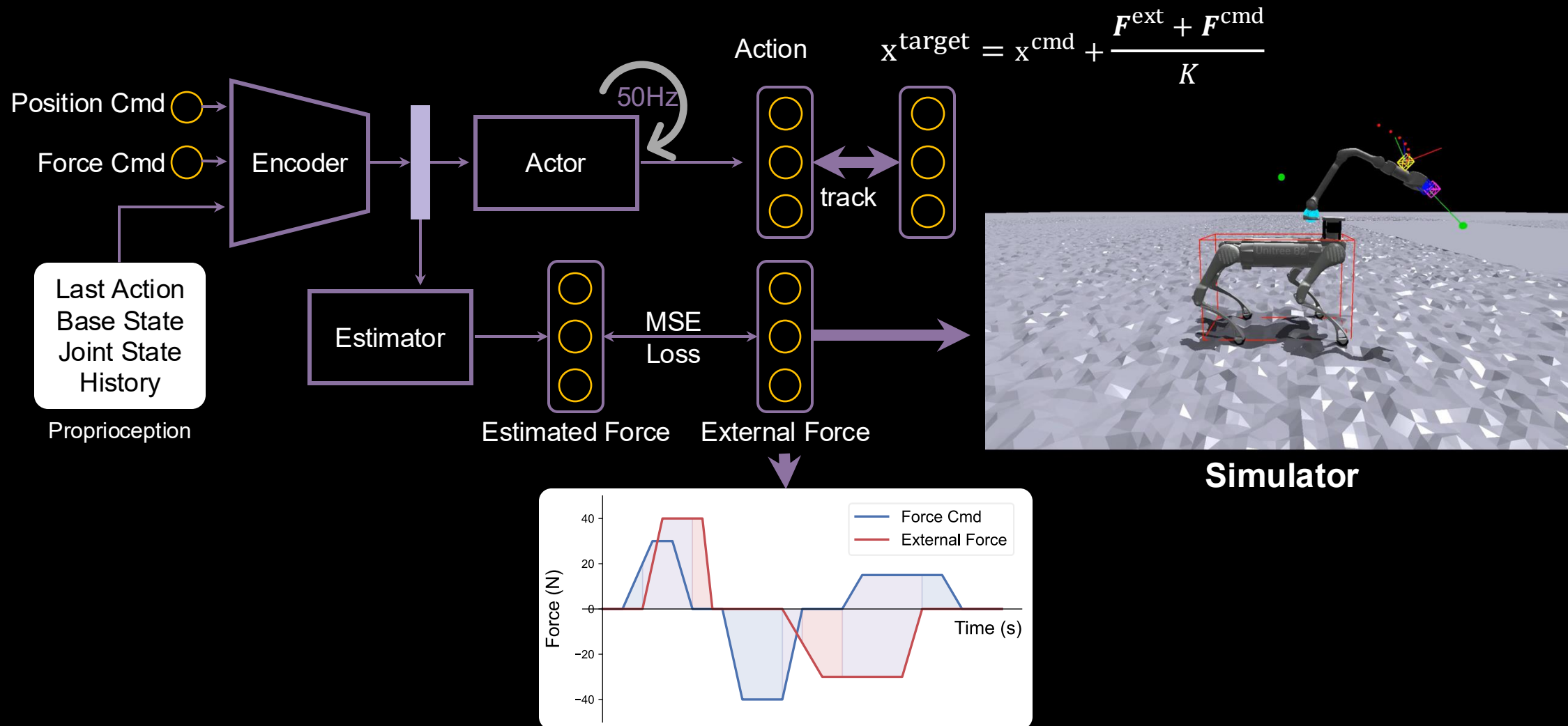
$$F = D(\dot{x} - \dot{x}^{\text{cmd}})$$

Compliant locomotion

$$\dot{x}^{\text{target}} = \dot{x}^{\text{cmd}} + \frac{F^{\text{ext}}}{D}$$

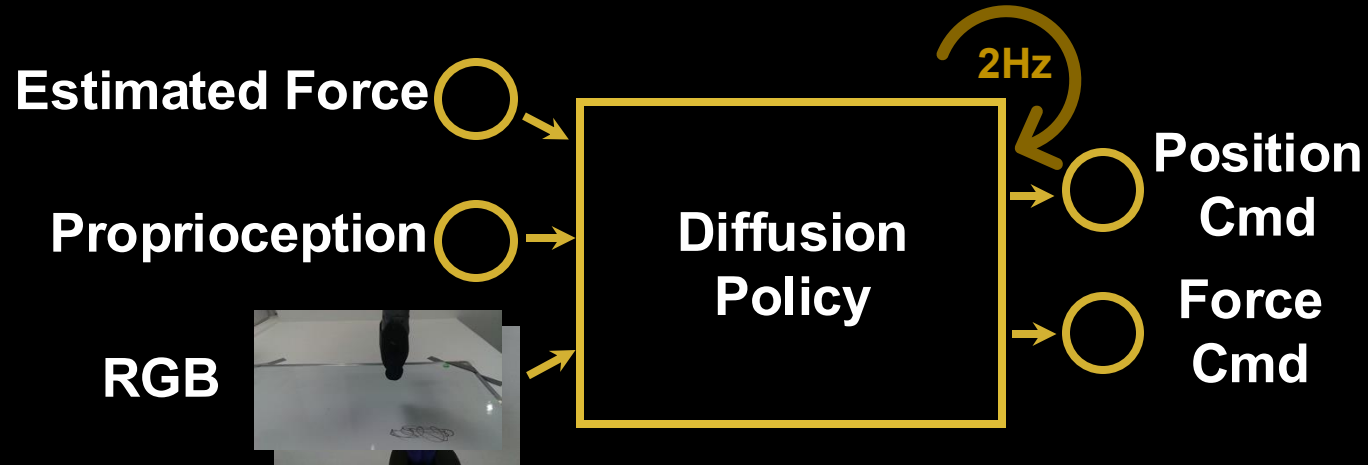


# UniFP via RL with force-position sampling in simulator

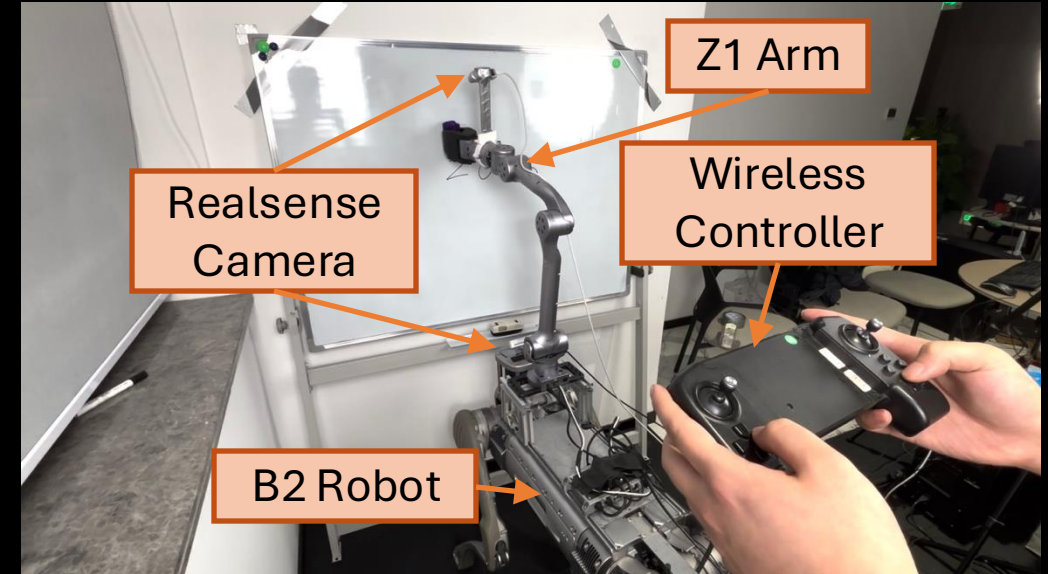




# UniFP for force-aware real-world imitation learning



- Data collection with **estimated forces**
- Imitation learning with **position and force command targets**
- Inference with **UniFP**



# UniFP for force-aware real-world imitation learning



**Tested on 4 tasks with each task taking 50 demonstrations**

# UniFP for force-aware real-world imitation learning

Table A.3: Imitation learning results (50 trials per task)

| Task      | wipe-blackboard | open-cabinet | close-cabinet | open-drawer-occlusion |
|-----------|-----------------|--------------|---------------|-----------------------|
| w/o Force | 0.22            | 0.36         | 0.30          | 0.30                  |
| w/ Force  | 0.58            | 0.70         | 0.72          | 0.76                  |



Base Camera View

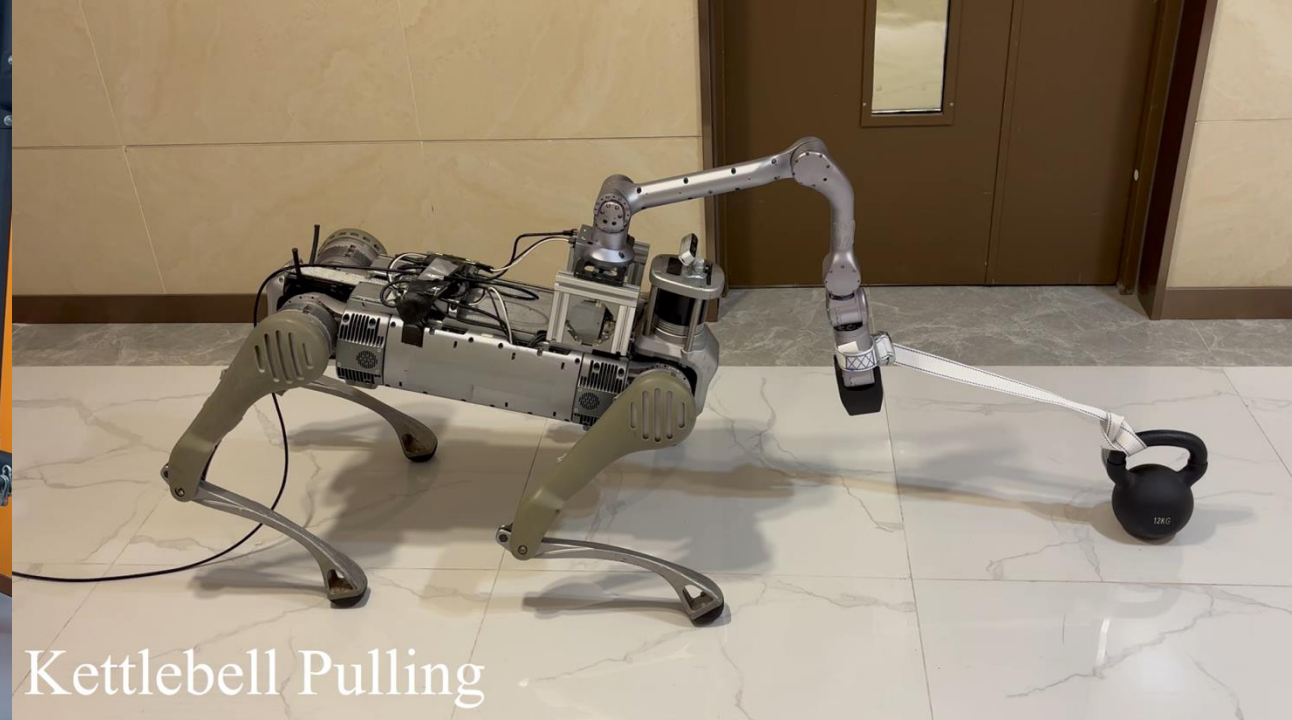


Achieves **~39.5%** higher success rate than the vanilla DP policy

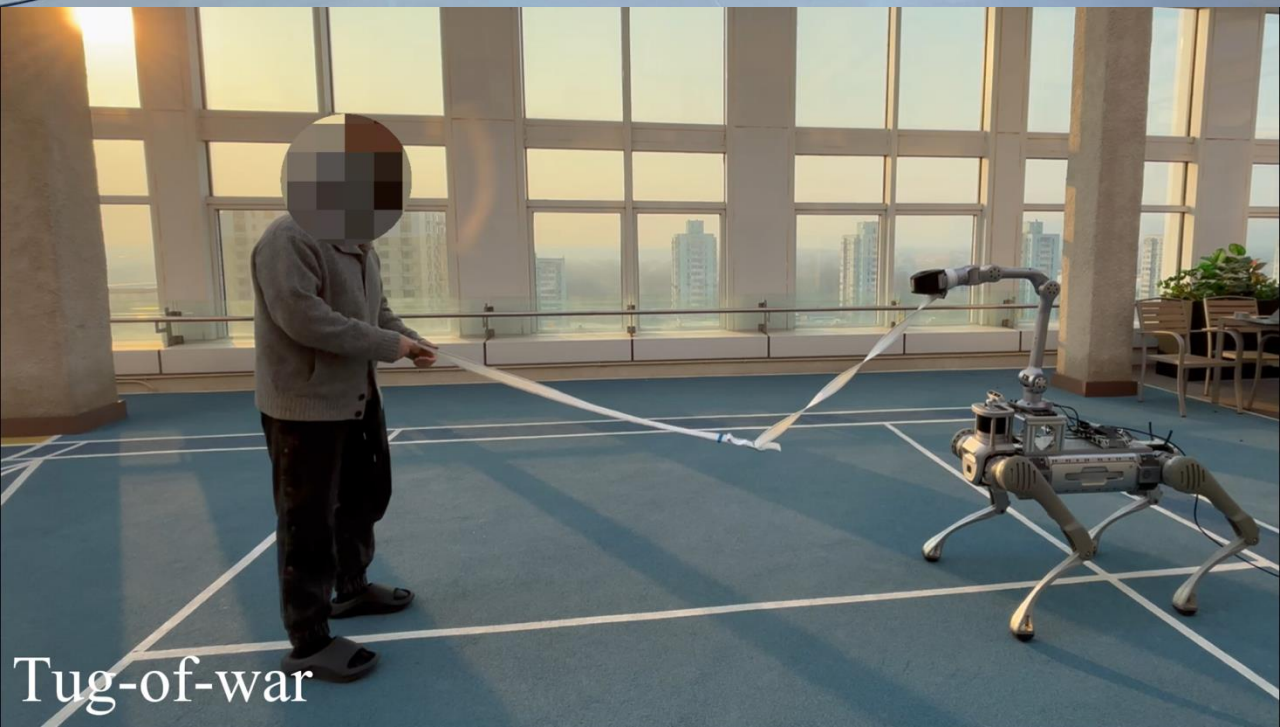




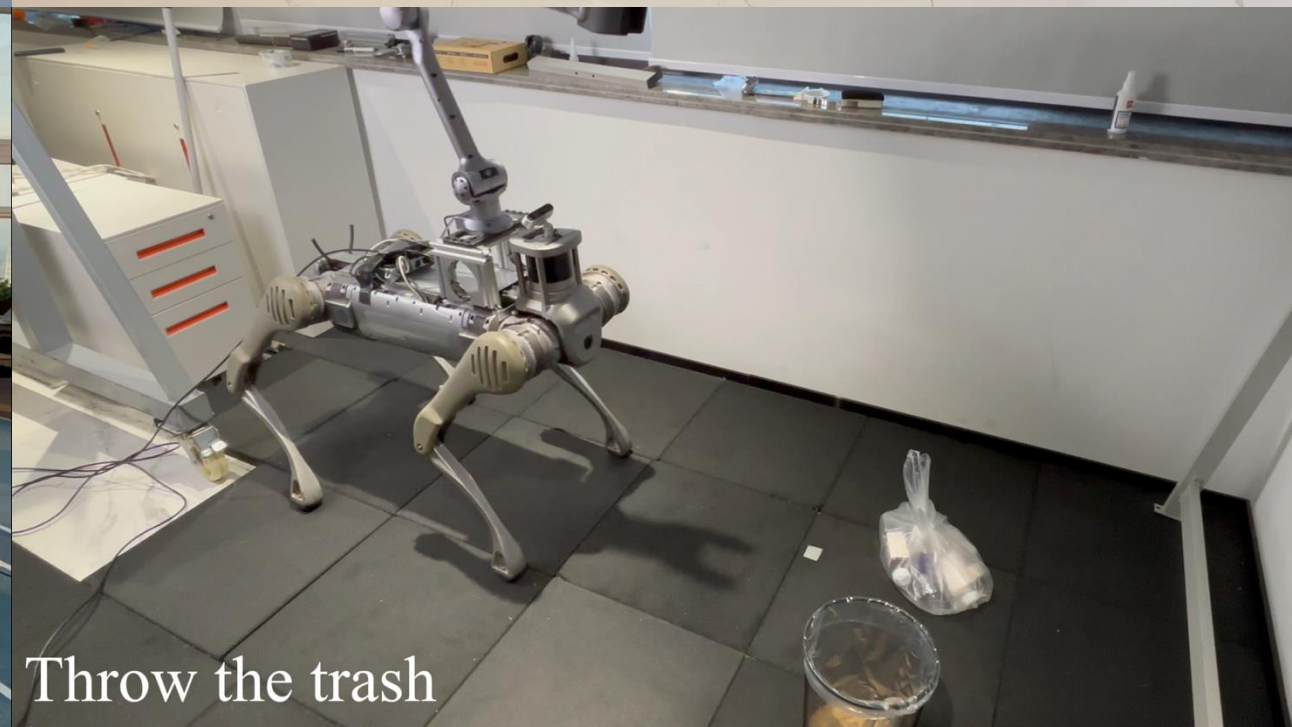
Robot Exercising in Gym



Kettlebell Pulling



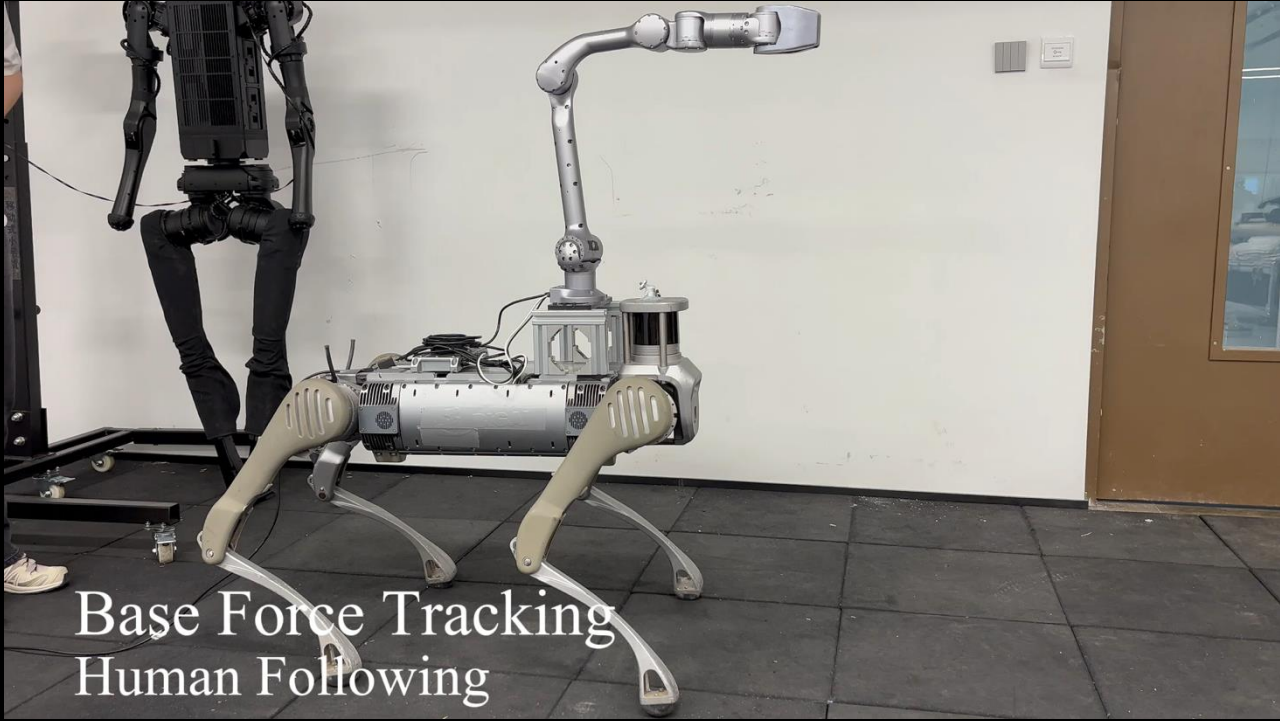
Tug-of-war



Throw the trash

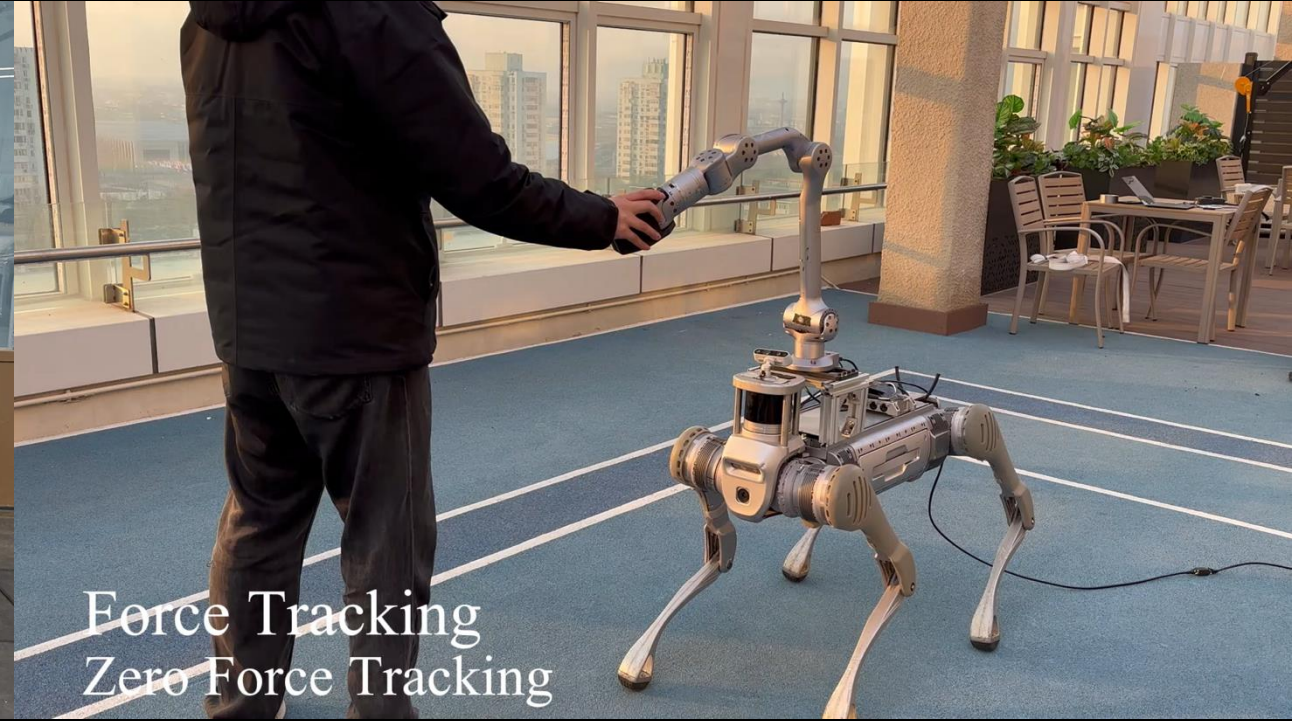


# So how is this important...



Base Force Tracking  
Human Following

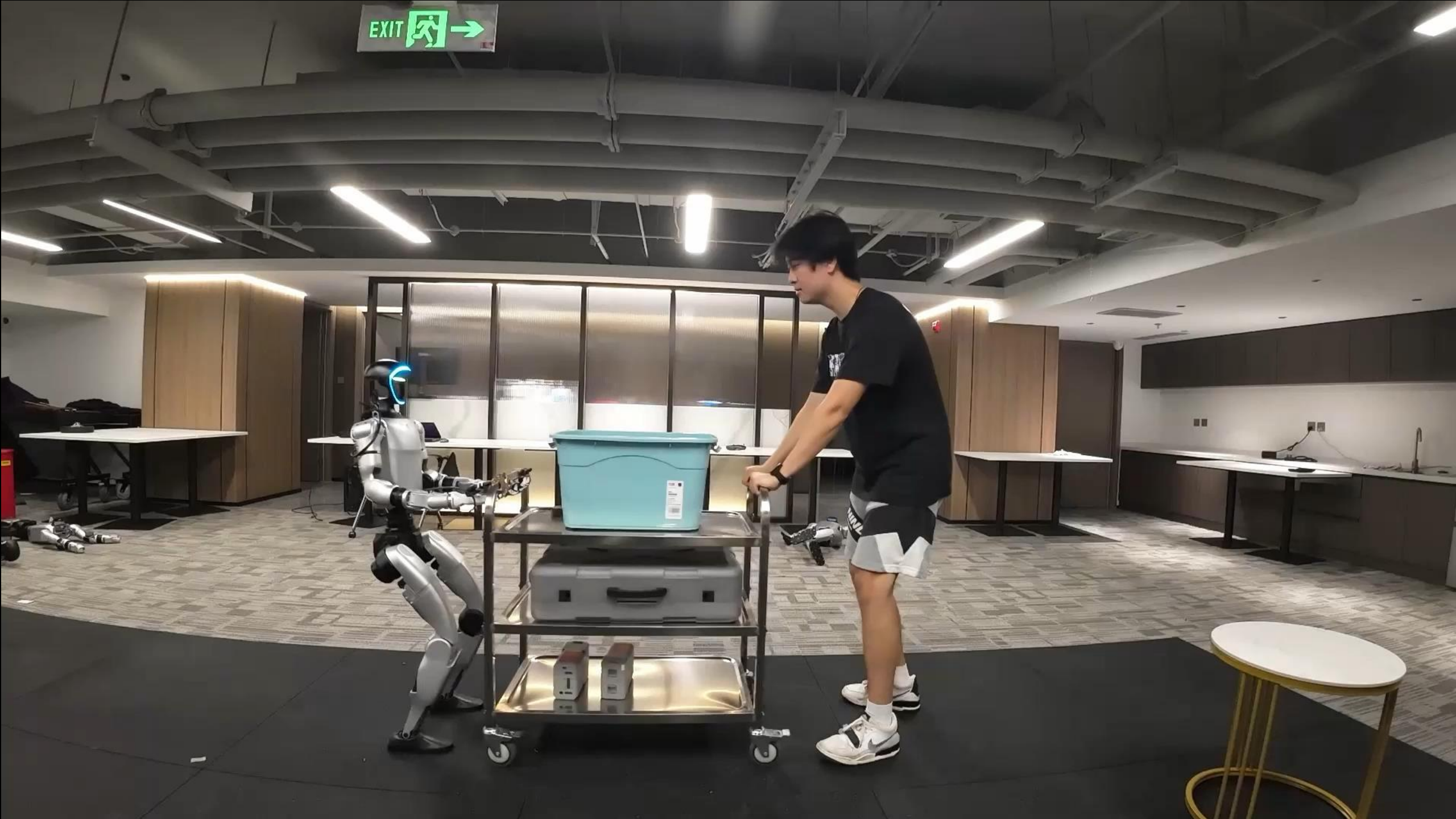
**Movement Tracking**



Force Tracking  
Zero Force Tracking

**Compliant Holding**

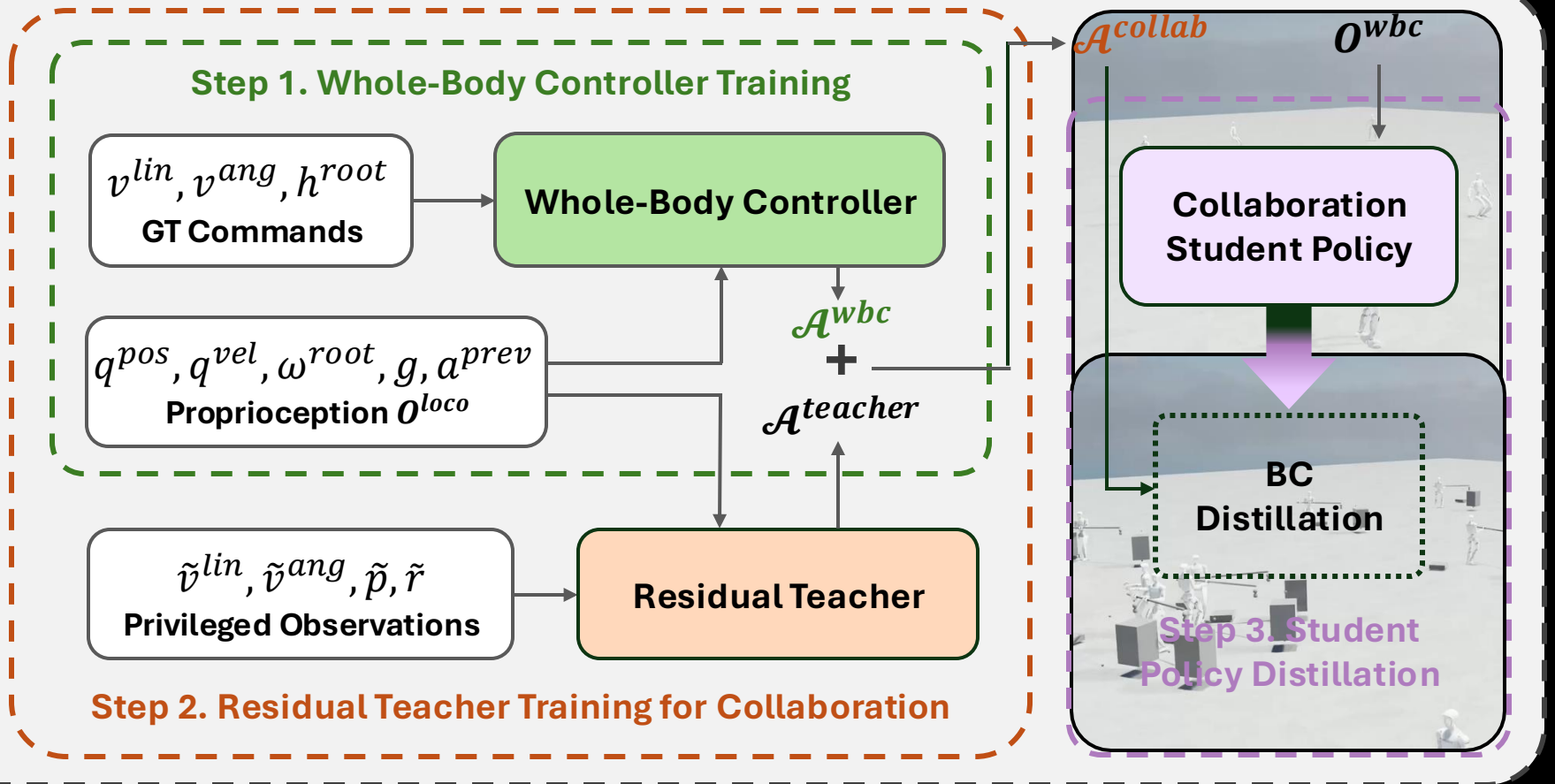
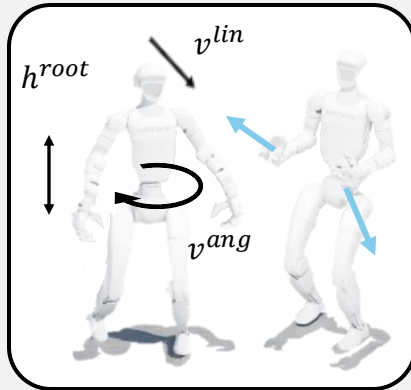
**The composition of these behaviors works for human-robot collaboration**



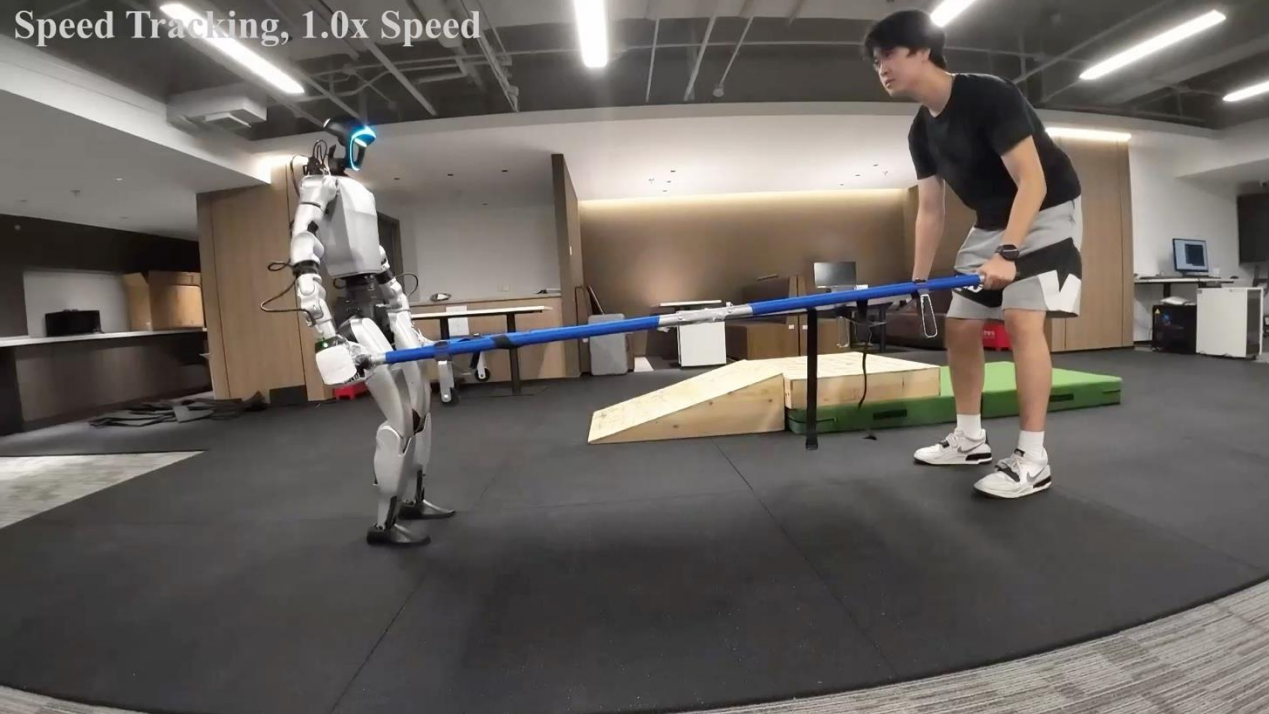


# COLA for collaborative object carrying

## External Forces



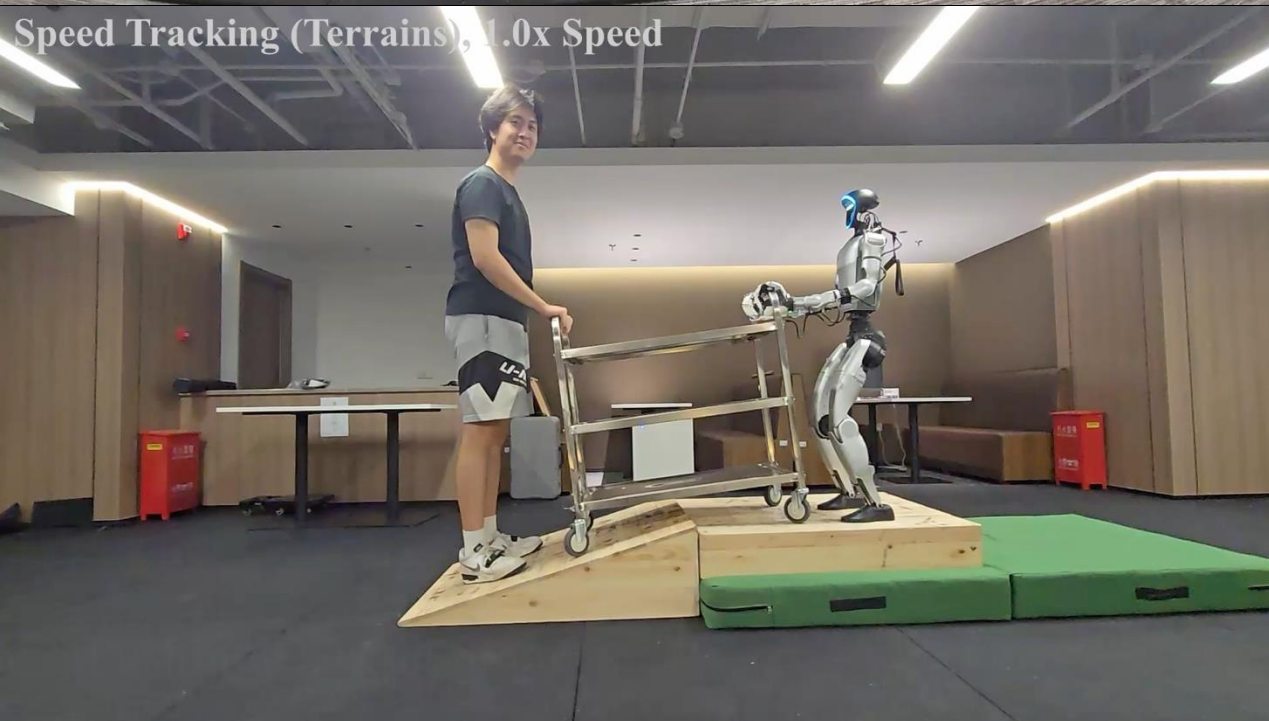
Speed Tracking, 1.0x Speed



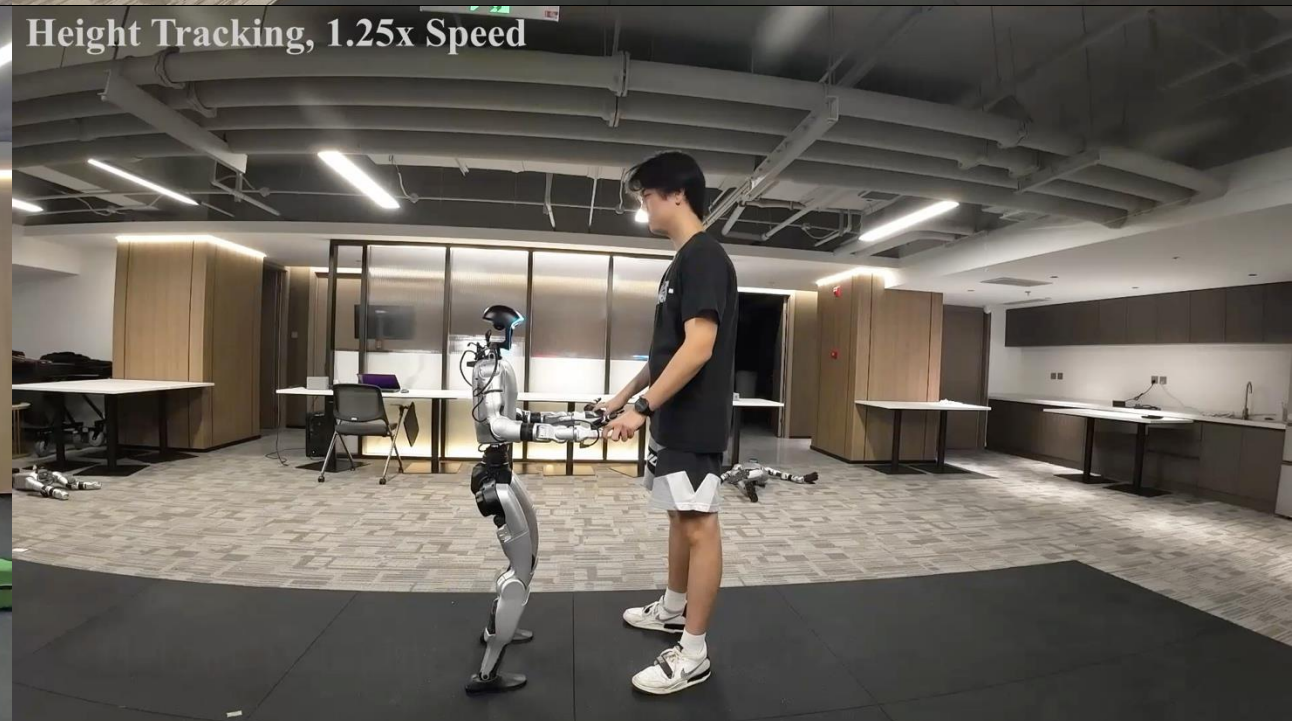
Height Tracking, 1.0x Speed



Speed Tracking (Terrains), 1.0x Speed



Height Tracking, 1.25x Speed





# Long Distance Testing (102.4m Total), 1.0x Speed







# Speed Tracking, 1.0x Speed

'Shopping Assistant' with Payload (15kg)

# Summary & Takeaways

- **Data scaling and unification as the main obstacle of EAI**
  - ❖ **Different embodiments, low-cost hardwares**
  - ❖ **Real2Sim2Real, synthetic augmentation, world models etc.**

# Summary & Takeaways

- **Data scaling and unification as the main obstacle of EAI**
  - ❖ Different embodiments, low-cost hardwares
  - ❖ Real2Sim2Real, synthetic augmentation, world models etc.
- **General reasoning and acting capabilities for robot tasks**
  - ❖ Aligning MLLMs for planning and interaction, efficient representations
  - ❖ Injecting spatial understanding capabilities for VLA models



# Summary & Takeaways

- **Data scaling and unification as the main obstacle of EAI**
  - ❖ Different embodiments, low-cost hardwares
  - ❖ Real2Sim2Real, synthetic augmentation, world models etc.
- **General reasoning and acting capabilities for robot tasks**
  - ❖ Aligning MLLMs for planning and interaction, efficient representations
  - ❖ Injecting spatial understanding capabilities for VLA models
- **Agile and safe robot control for human-robot interaction**
  - ❖ Recover the missing force modality for compliance policies
  - ❖ Safe control behaviors over VLA for human-robot interaction

RoboVerse (RSS 2025)

DP-Recon (CVPR 2025)

"Pokemon style"

Diverse Tasks and Demonstrations

"Van Gogh style"

"Minecraft style"

RoboVerse

Thank you  
Q&A

ManipTrans (CVPR 2025)

ControlVLA (CoRL 2025)

Webpage

WeChat

