# Generate Robotic Data with Spatial Intelligence
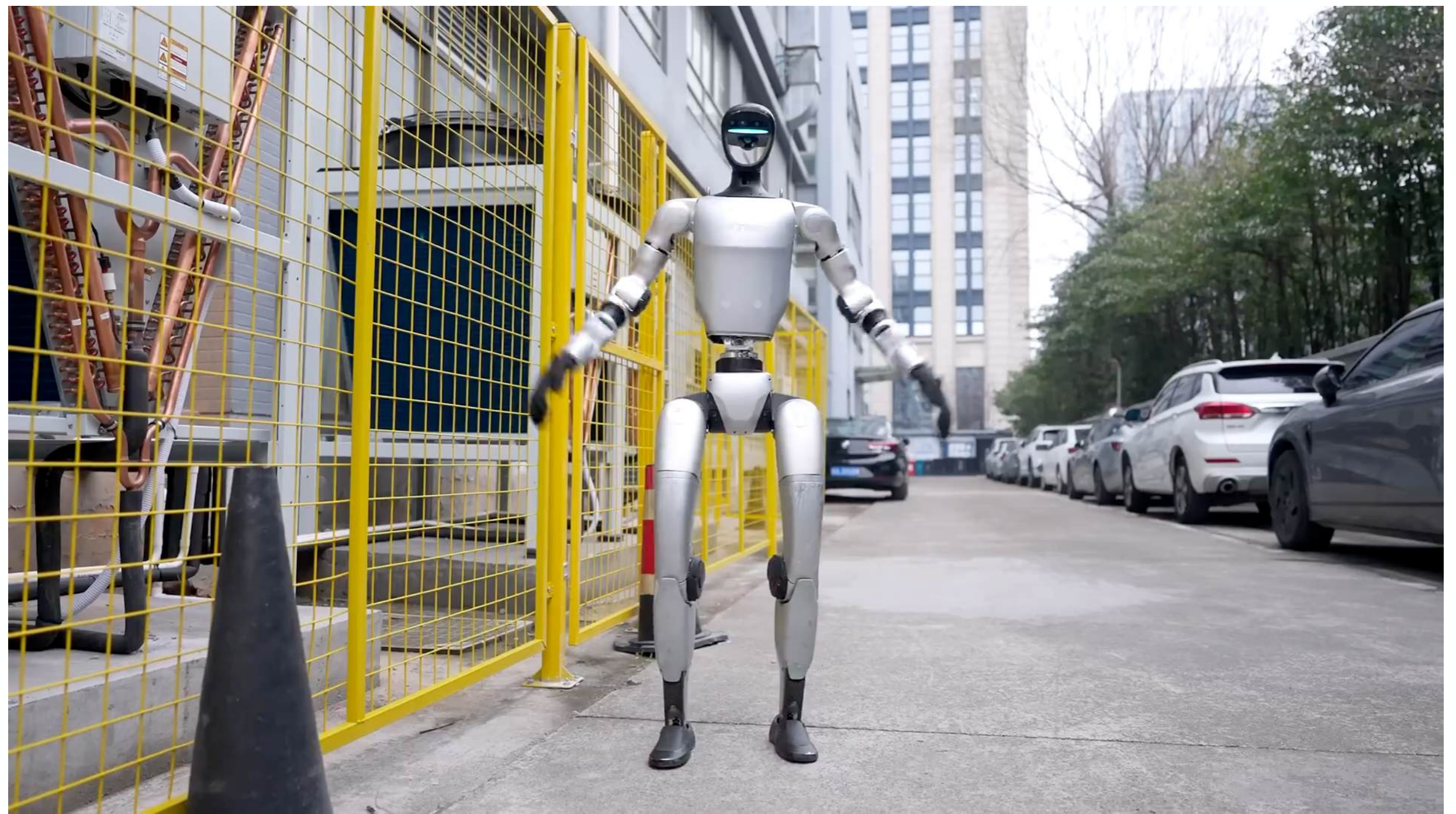
Yue Wang
MUSI | Oct 20th, 2025
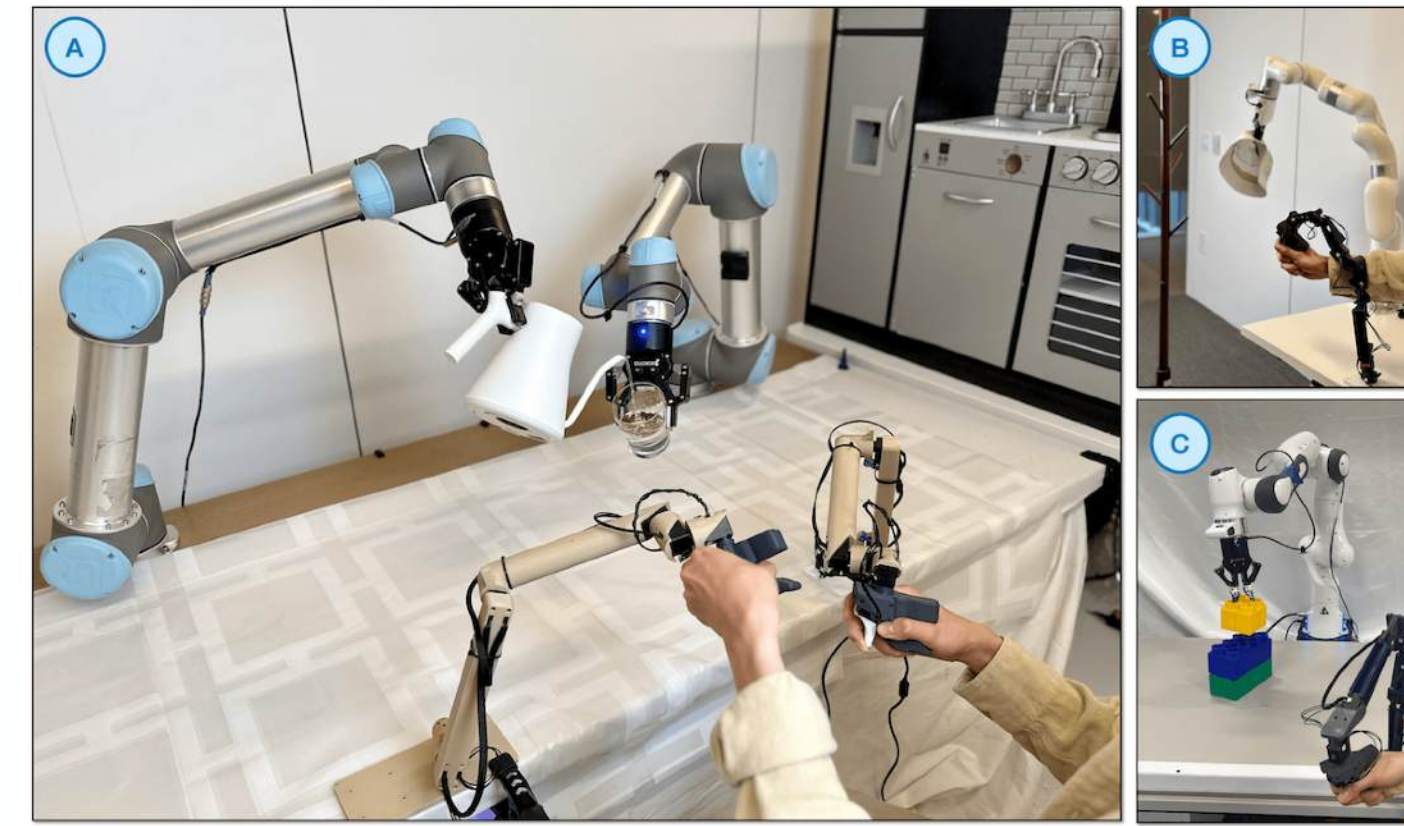
# Cambrian Explosion of Robotics

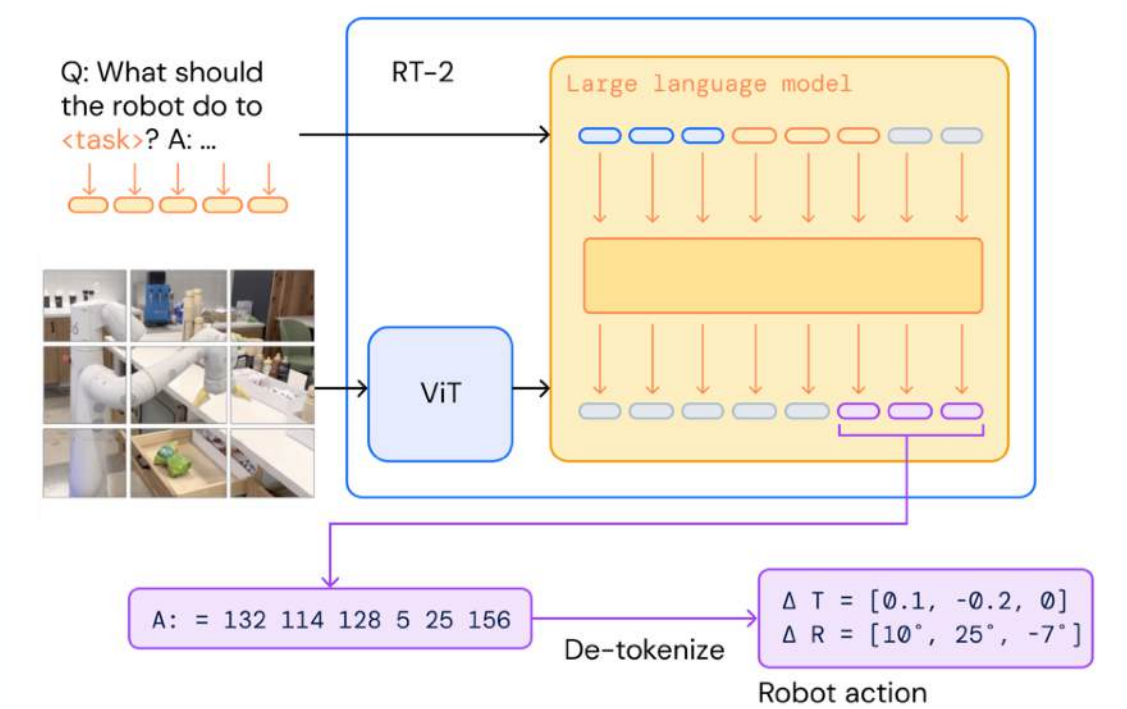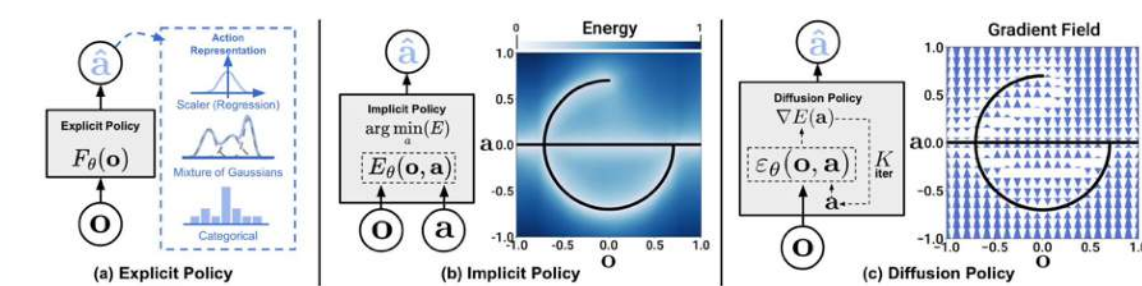1x speed, autonomous
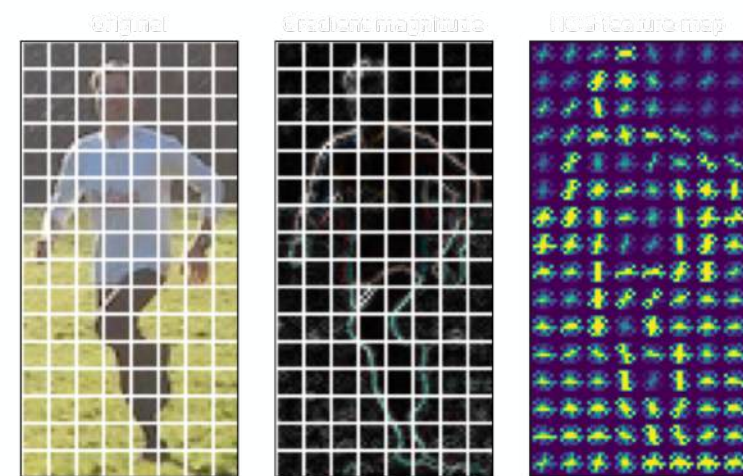
Generalist

2

Data

Hardware

Algorithm

Diffusion Policy
Visuomotor Policy Learning via Action Diffusion
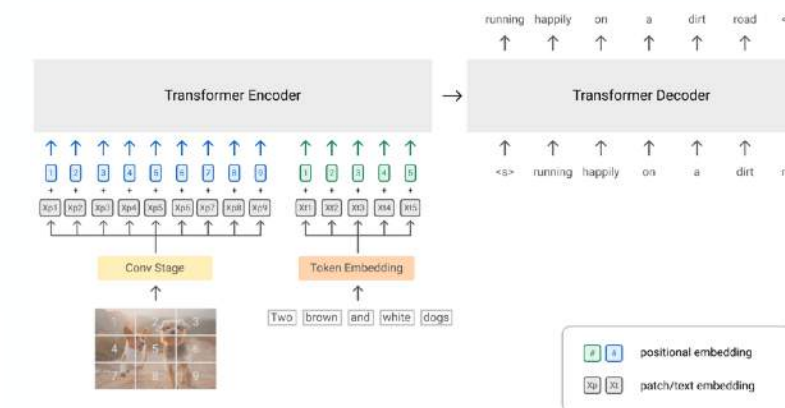
Data is the key to artificial intelligence.

Physical AI

VLMs

Transformers

CNNs

HOG+SIFT+SVM

Little Data    Curated    Web Scale    Multimodal    Embodied Data

| < 1s | > 60s |
| --- | --- |
| Ubiquitous | Confined to lab environments |
| $0.01 per data point | $5 per data point |

How to generate robotic data with spatial intelligence techniques?

# How to generate robotic data with spatial intelligence techniques?

## Use Real-to-Sim Reconstruction



Robot Learning from Any Images. Zhao et al. CoRL 2025.

## Lever[age]



Learning from Massive Human Videos for Universal Humanoid Pose Control. Mao et al. Humanoids 2025.

## Scale Teleoperation Data



Humanoid Everyday. Jing et al. In submission.

# Robot Learning from Any Images

Robot Learning from Any Images. Zhao et al. CoRL 2025.

# Step-1: Recovering the Physical Scene from a Single Image



Input Image $I$

On-device Camera

Mobile Phone Capture

Robotic Dataset

Internet Image

# Step-1: Recovering the Physical Scene from a Single Image



Input Image $I$

Segmentation & Inpainting

Metric Depth Prediction
& Point Cloud Generation

Step-1: Recovering the Physical Scene from a Single Image

Input Image $I$

Segmentation & Inpainting

Metric Depth Prediction & Point Cloud Generation

Geometry & Appearance Modeling

Step-1: Recovering the Physical Scene from a Single Image

Input Image $I$

Segmentation & Inpainting

Metric Depth Prediction & Point Cloud Generation

Geometry & Appearance Modeling

Recovering Scene Configuration

Step-1: Recovering the Physical Scene from a Single Image

Input Image $I$

Segmentation & Inpainting

Metric Depth Prediction & Point Cloud Generation

Geometry & Appearance Modeling

Recovering Scene Configuration

Physical Property Estimation & Robot Placement

Step-1: Recovering the Physical Scene from a Single Image

Input Image $I$

Segmentation & Inpainting

Metric Depth Prediction & Point Cloud Generation

Geometry & Appearance Modeling

Recovering Scene Configuration

Physical Property Estimation & Robot Placement

Step-2: Scalable Robotic Data Generation in Sim

Robotic Data Generation

Step-1: Recovering the Physical Scene from a Single Image

Input Image $I$

Segmentation & Inpainting

Metric Depth Prediction & Point Cloud Generation

Geometry & Appearance Modeling

Recovering Scene Configuration

Physical Property Estimation & Robot Placement

Step-2: Scalable Robotic Data Generation in Sim

Visual Blending

Robotic Data Generation

14

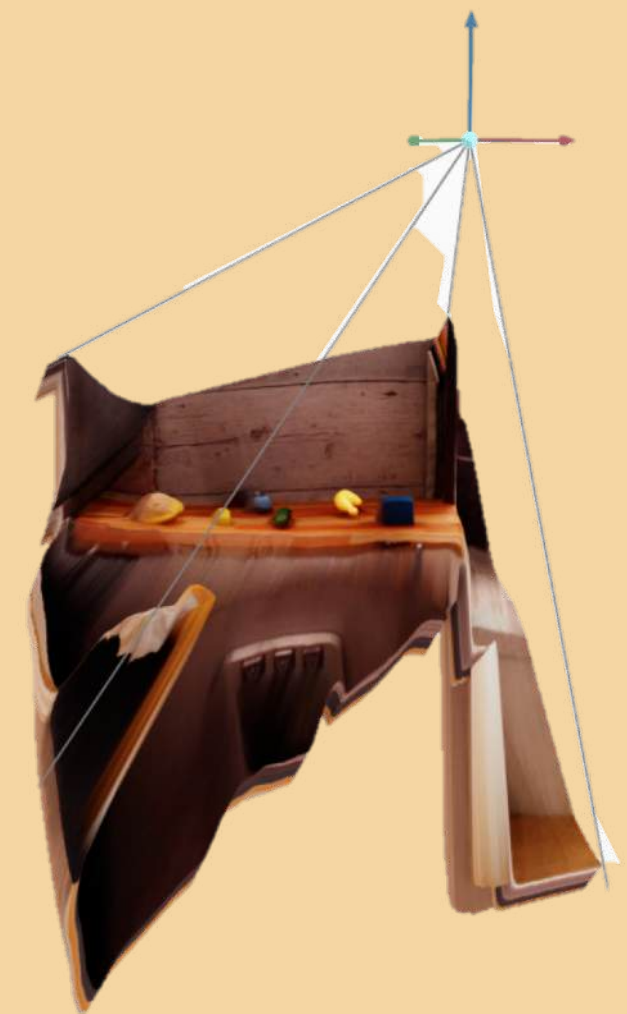**Step-1: Recovering the Physical Scene from a Single Image**

Input Image $I$

Segmentation & Inpainting

Metric Depth Prediction & Point Cloud Generation

Geometry & Appearance Modeling

Recovering Scene Configuration

Physical Property Estimation & Robot Placement

**Step-3: Robot Learning & Deployment**

Real-World Deployment

Single-Image IL

VLA

Manip. Priors

Robotic Data

Robotic Images

Camera Photographs

Internet Images

**Step-2: Scalable Robotic Data Generation in Sim**

Visual Blending

Robotic Data Generation

# Single-Image Imitation

RoLA-Generated Data @ Sim

Real-world Deploy

## Manipulation in Cluttered Scenes



## Pour Water

# Data Collection

# Real-World Deployment



Pick up the carrot.

Pick up the yellow banana.

Put the yellow lemon beside the green apple.

Take the grey object beside the lemon and place it beside the yellow apple.

# Manipulation Prior 🧠

# Robot Learning from Any Images

🗄 Data quantity and diversity are widely recognized as primary bottlenecks in scaling robot learning.

🦾 Collecting **on-robot demonstrations** at scale demands specialized hardware and extensive labor. ☹

🖼 Obtain robot-complete data from non-robotic images under minimal assumptions: **single image**. ☺

# Robot Learning from A Physical World Model



Image & Task Prompt

Task: Pour the tomato from the pan into the white plate

Task-Conditioned Video Generation

Physical World Modeling & Learning

4D Point Cloud Reconstruction    Textured Mesh Generation    Physical Scene Reconstruction

Object-Centric Learning from the Physical World Model

Zero-Shot Manipulation in the Real World

# Robot Learning from A Physical World Model



Video generation



Robot execution

Wrist: 3 DoFs

Elbow: 1 DoF

Shoulder: 3 DoFs

Leg: 6 DoFs

More DoFs
Not easily handled by motion model
Action retargeting is hard

How can we derive humanoid data from Internet data?

# UH-1: Learning from Massive Human Videos for Universal Humanoid Pose Control

**Massive Internet Videos**

YouTube

Youtube

DeepMind

Kinetics 700

AI2 Allen Institute for AI

Charades

**Video Clip Extraction**



**3D Human Pose Estimation**



**Video Captioning**



*"A young woman is doing a workout in a living room, using her legs and arms to perform various exercises."*

**Goal-based Reinforcement Learning**

**Motion Retargeting from Humans to Humanoids**

**Text Representation**



*"A young woman is doing a workout in a living room, using her legs and arms to perform various exercises."*

**Humanoid Robot Actions**

**UH-1 Transformer**

**Humanoid Action Tokens**

An automatic humanoid data engine

A unified whole-body control model

"Learning from Massive Human Videos for Universal Humanoid Pose Control." Mao et al. Humanoids 2025.

# Data Collection

We collect 163, 800 video clips from diverse sources.

# Data Collection

## Videos are further annotated with captioning tools.



"practicing martial arts, standing."

"standing and speaking."

"practicing yoga."

"doing squats, lunges, and jumping jacks."

**VideoLLaMA 2:** The video features *a kitten and a baby chick* playing together. They are seen *cuddling, playing, and even taking a nap* together. The video has a very *cute and heartwarming* feel to it, as the two animals seem to have *formed a close bond*.

Pre-trained Large Language Model

**Prompt:** What animals are in the video, what are they doing, and how does the video feel?

Projection $W$

Projection $W$

Flatten

Flatten

Spatial Convolution

Audio Encoder

Spatial-Temporal Downsampling

Spatial Convolution

Visual Encoder

Video Frames          Encoding          STC connector          Audio

[Preprint 2024] "VideoLLaMA 2: Advancing Spatial-Temporal Modeling and Audio Understanding in Video-LLMs." Cheng et al.

# Human Motion Representation



Shape: PCA coefficients

$\vec{\beta}$

Pose: Rotation of 23 joints

$\vec{\theta}$

SMPL Model

HMR 2.0

Vision Transformer

Input Image

SMPL Query Token

Transformer w/ Cross Attn

MLP

$\theta$ Pose
$\beta$ Shape
$\pi$ Camera

Generator ($\mathcal{G}$)

CNN

GRU

time

[Kocabas et al., CVPR 2020]
[Goel et al., CVPR 2024]

# Human-to-Humanoid Motion Retargeting



$$\min_{\beta} ||\mathcal{P}^T_{joints} - \mathcal{P}^T_{robot}||_2,$$

$$\text{s.t.} \quad \mathcal{P}^T_{joints} = F_{fk}(\mathcal{P}_{human}(\beta, \theta^T, t_{root})),$$

# Sim-to-Real Adaptation



PD Controller *(200Hz)*

| Simulation | Privileged Proprioception 361 dim | Privileged Imitation Policy $\pi_{\text{privileged}}$ *50Hz* | $a_t^{\text{privileged}}$ |
| | Privileged Motion Goal 552 dim | | |

*Phase 1* *Reinforcement Learning*

*Phase 2* *Supervised Learning*

| Imitation Goal | Sim-to-Real Proprioception 26 history step *63 dim=1638 dim | Sim-to-Real Imitation Policy $\pi_{\text{OmniH2O}}$ *50Hz* | $a_t$ |
| | Sim-to-Real Motion Goal 27 dim | | |

[He et al., CoRL 2024]

# Dataset



(a) Data Distribution

(b) Data Scale

(c) Vocabulary Diversity

(d) Motion Sample

# Universal Humanoid (UH-1) Architecture

**Text Representation**

*"A young woman is doing a workout in a living room, using her legs and arms to perform various exercises."*

**UH-1 Transformer**

**Humanoid Robot Actions**

**Humanoid Action Tokens**

**Encoder** → $\mathcal{Z}$ → **Decoder**

**UH-1 Humanoid Action Tokenizer**

**User**

*Wave hand*
*Open bottle*
*Play violin*
*...*

**UH1**

Text-to-Keypoint

Text-to-Action

RL Policy

Execution

One Model, Two Control Modes

# Research Questions

- **Universal Pose Control with UH-1**: Does UH-1 model enable universal humanoid robot pose control based on text commands?

- **Scalability and Generalization with Humanoid-X**: Does the large-scale Humanoid-X dataset facilitate scalable training and improve the generalization ability of UH-1?

- **Real-World Deployment of UH-1**: Can UH-1 model be deployed on real humanoid robots to enable reliable robotic control in real-world environments?

# Universal Pose Control with UH-1

- Baseline models: Motion Diffusion Model (MDM) and Text-to-Motion GPT (T2M-GPT)

| Methods | FID ↓ | MM Dist ↓ | Diversity ↑ | R Precision ↑ |
|---|---|---|---|---|
| Oracle | $0.005^{\pm.001}$ | $3.140^{\pm.010}$ | $9.846^{\pm.062}$ | $0.780^{\pm.003}$ |
| MDM [57] | $0.582^{\pm.051}$ | $5.921^{\pm.034}$ | $10.122^{\pm.078}$ | $0.617^{\pm.007}$ |
| T2M-GPT [71] | $0.667^{\pm.109}$ | $3.401^{\pm.017}$ | $\mathbf{10.328^{\pm.099}}$ | $0.734^{\pm.004}$ |
| UH-1 (ours) | $\mathbf{0.445^{\pm.078}}$ | $\mathbf{3.249^{\pm.016}}$ | $10.157^{\pm.106}$ | $\mathbf{0.761^{\pm.003}}$ |

# Scalable Learning with Humanoid-X

- Increasing data size leads to consistent performance improvement.

- Pre-training on Humanoid-X helps generalization.



(a) FID ↓



(b) Diversity ↑

| Dataset | FID ↓ | MM Dist ↓ | Diversity ↑ | R Precision ↑ |
|---------|-------|-----------|-------------|---------------|
| Oracle | $0.005^{\pm.001}$ | $3.140^{\pm.010}$ | $9.846^{\pm.062}$ | $0.780^{\pm.003}$ |
| HumanoidML3D | $0.445^{\pm.078}$ | $3.249^{\pm.016}$ | $10.157^{\pm.106}$ | $0.760^{\pm.003}$ |
| Humanoid-X | $\mathbf{0.379^{\pm.046}}$ | $\mathbf{3.232^{\pm.008}}$ | $\mathbf{10.221^{\pm.100}}$ | $\mathbf{0.761^{\pm.003}}$ |

# Scalable Learning with Humanoid-X

- Increasing data size leads to consistent performance improvement.

- Pre-training on Humanoid-X helps generalization.



(a) FID ↓
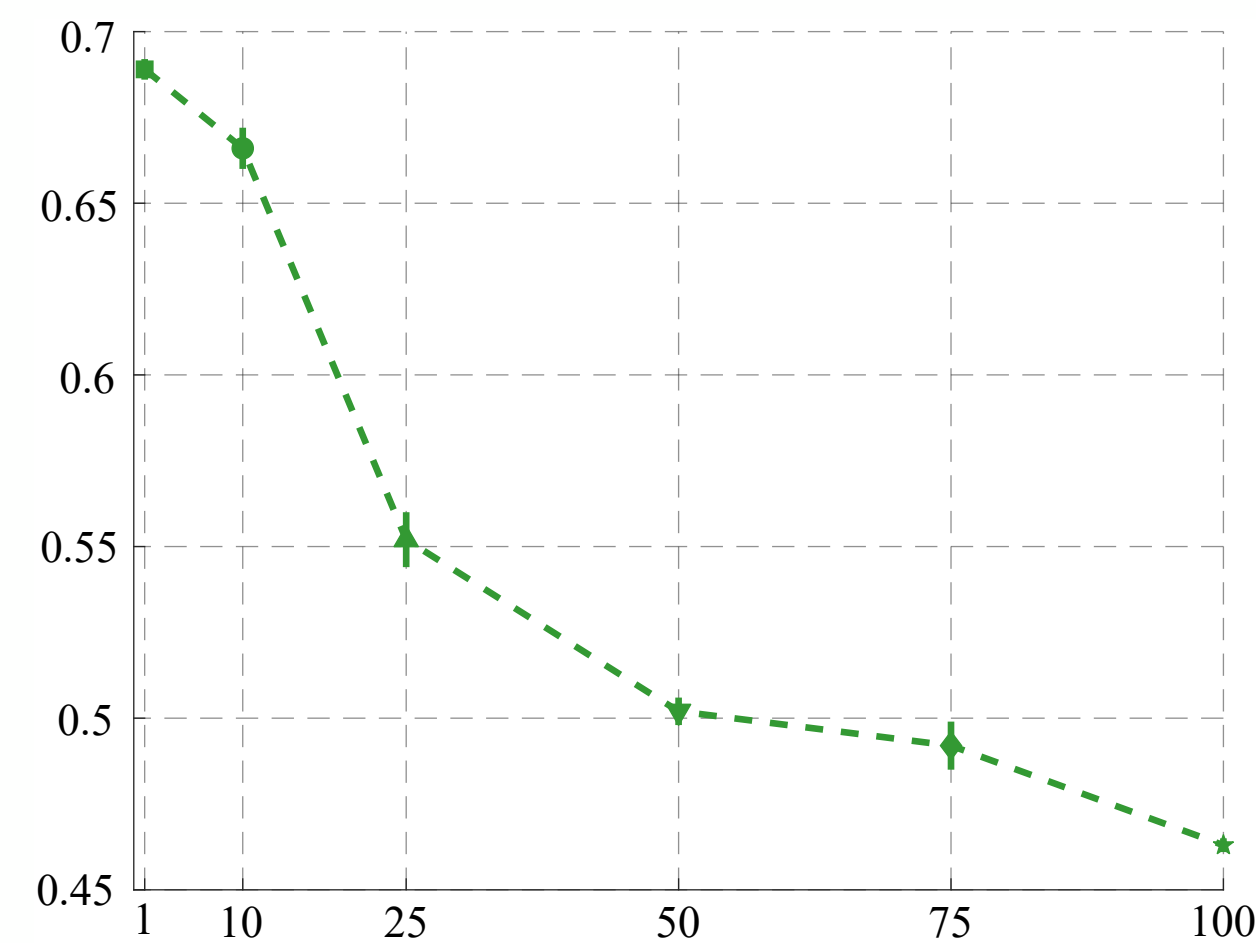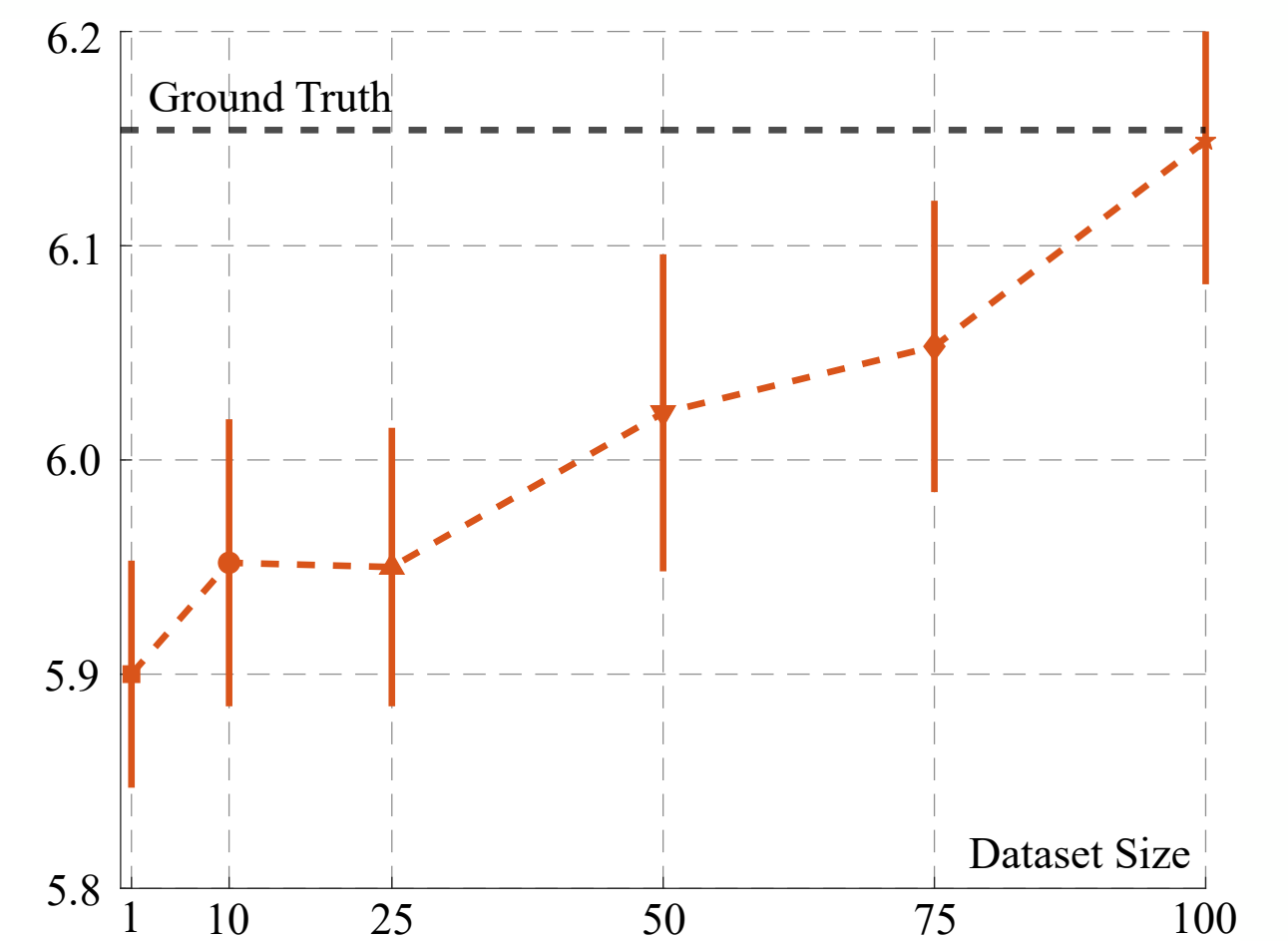


(b) Diversity ↑

| Dataset | FID ↓ | MM Dist ↓ | Diversity ↑ | R Precision ↑ |
|---|---|---|---|---|
| Oracle | $0.005^{\pm.001}$ | $3.140^{\pm.010}$ | $9.846^{\pm.062}$ | $0.780^{\pm.003}$ |
| HumanoidML3D | $0.445^{\pm.078}$ | $3.249^{\pm.016}$ | $10.157^{\pm.106}$ | $0.760^{\pm.003}$ |
| Humanoid-X | $\mathbf{0.379^{\pm.046}}$ | $\mathbf{3.232^{\pm.008}}$ | $\mathbf{10.221^{\pm.100}}$ | $\mathbf{0.761^{\pm.003}}$ |

# Real-World Deployment of UH-1

# Humanoid
# Everyday

A high-frequency humanoid dataset spanning diverse everyday tasks.
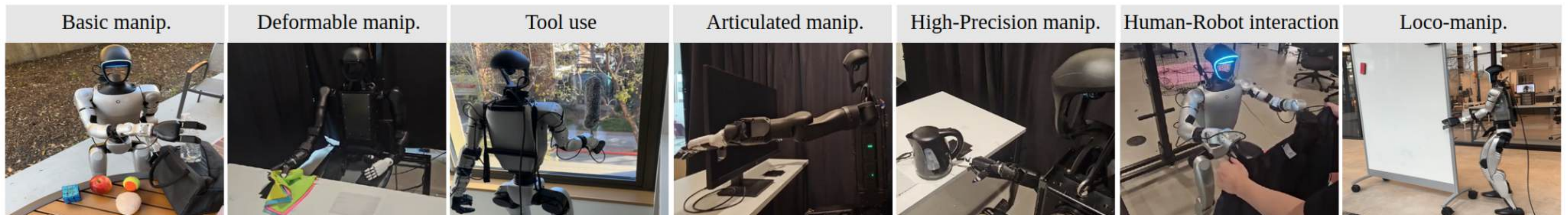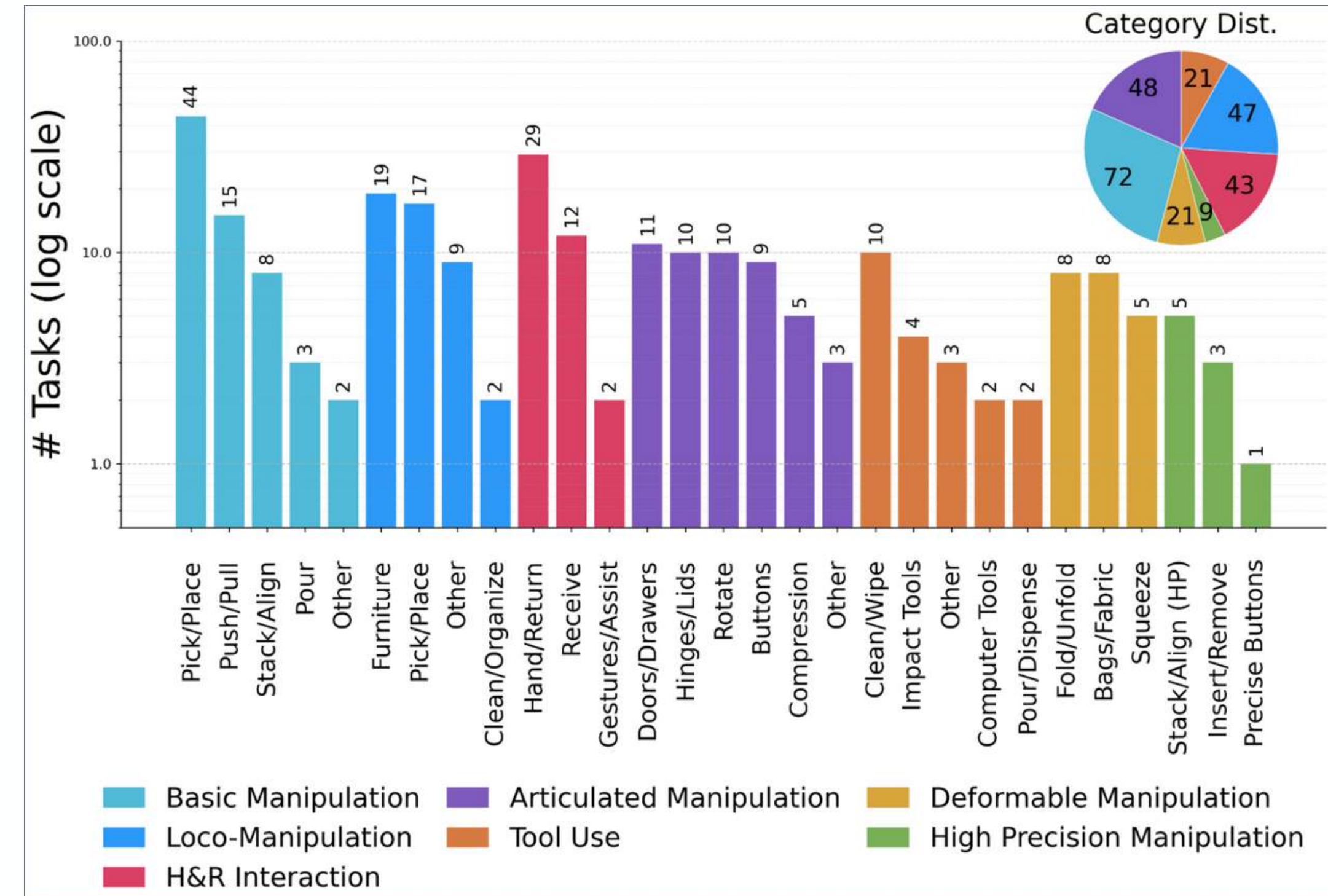
- 260 Tasks
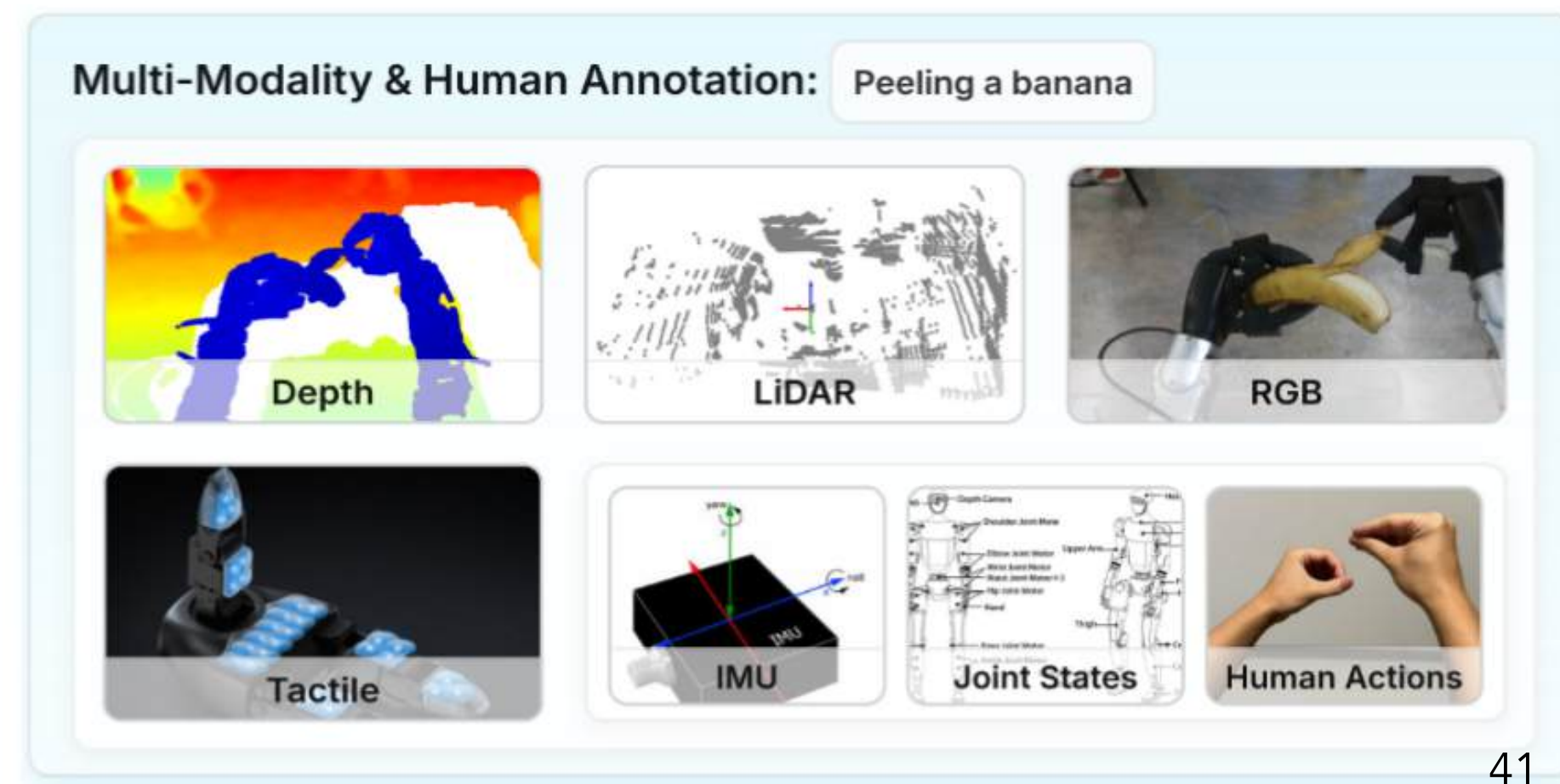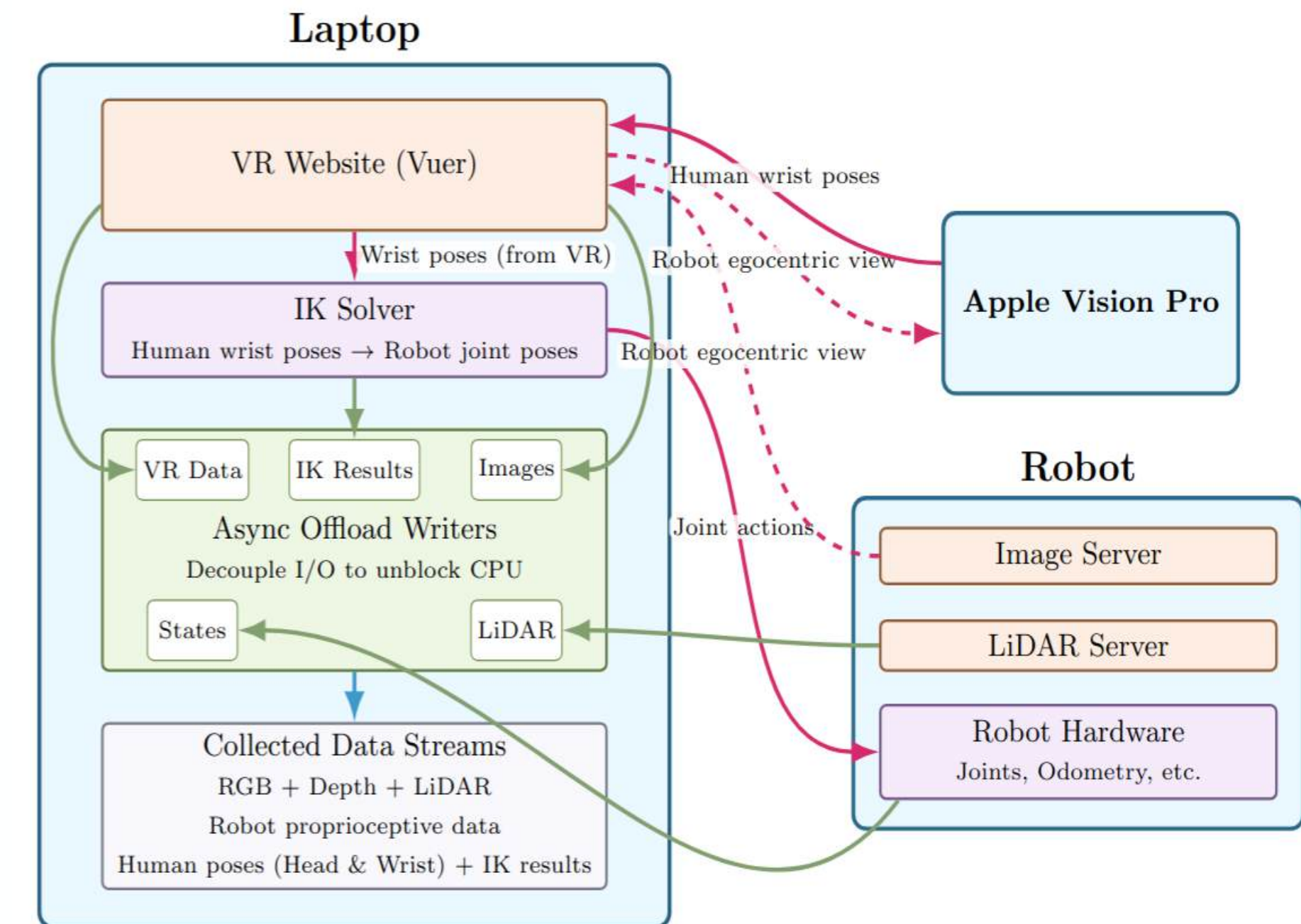- 7 categories
- 30 Hz
- 9 Modalities
- Open Source

# Dataset: Diverse collection of humanoid tasks

- Covers 10.3K trajectories, 3M+ frames, and 260 tasks using Unitree G1 and H1

- Includes bipedal loco-manipulation and human-robot Interaction that are rare in other datasets
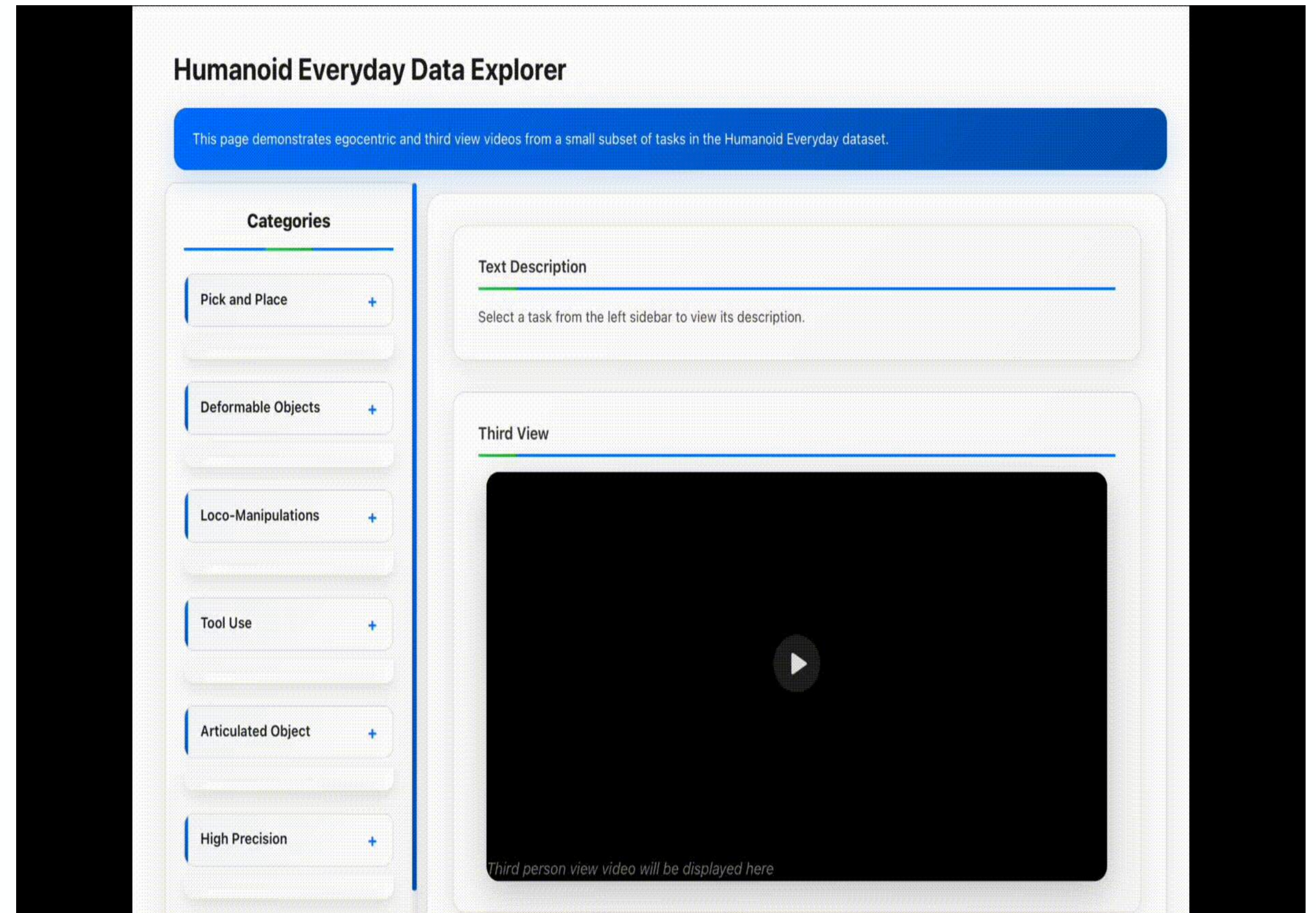
# Dataset: efficient data collection pipeline

- ## Collection Pipeline:
  - Offloaded I/O keeps control loop fast and responsive.

- ## Improved Performance:
  - Reduced control delay from 500ms to 20ms
  - Halved data collection time

- ## 30hz multi-modality streams collected:
  - RGB+Depth+LiDAR
  - Proprioceptives: Joint States, Tactile, Odometry, IMU
  - Human Actions+Task Descriptions



Laptop

VR Website (Vuer)

Human wrist poses

Wrist poses (from VR)    Robot egocentric view

Apple Vision Pro

IK Solver
Human wrist poses → Robot joint poses    Robot egocentric view

VR Data    IK Results    Images
Async Offload Writers
Decouple I/O to unblock CPU

States    LiDAR

Joint actions

Robot

Image Server

LiDAR Server

Robot Hardware
Joints, Odometry, etc.

Collected Data Streams
RGB + Depth + LiDAR
Robot proprioceptive data
Human poses (Head & Wrist) + IK results



Multi-Modality & Human Annotation: Peeling a banana

Depth    LiDAR    RGB

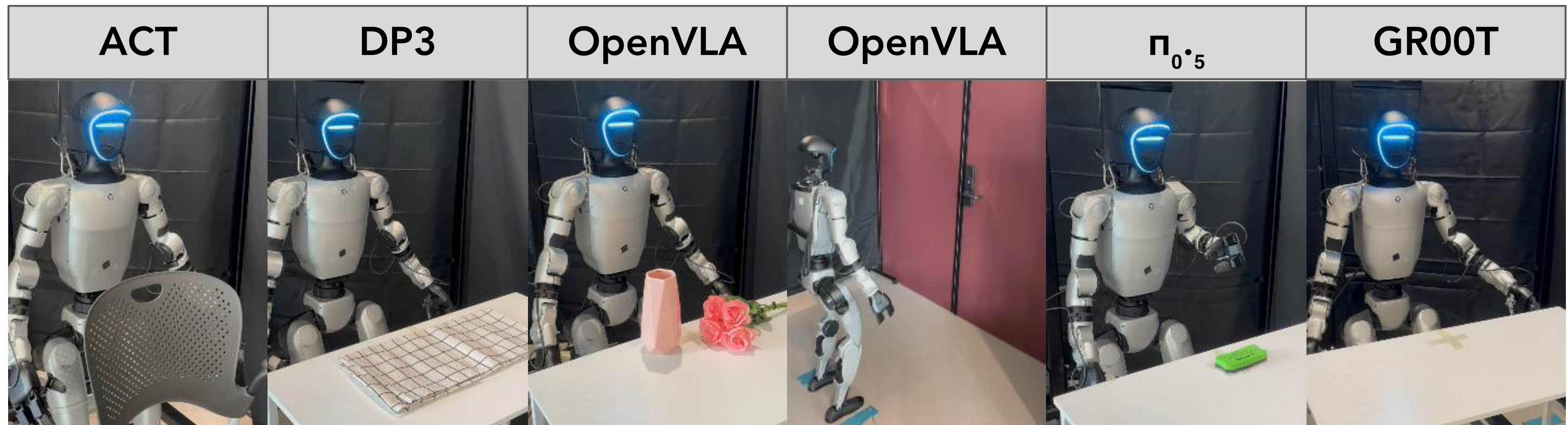Tactile    IMU    Joint States    Human Actions

# Dataset: Data viewer

- Data viewer contains 50 sample tasks from all of our categories
- Structure
  - Text Description
  - Third View Video
  - Egocentric Video
  - Point Cloud/Depth Visualization

# Policy Inference: Imitation Learning + VLA

- We run inference using different imitation learning policies and VLA models on different manipulation tasks.



| ACT | DP3 | OpenVLA | OpenVLA | $\pi_0 \cdot_5$ | GR00T |
|-----|-----|---------|---------|----------|-------|

# Results

| Task Category | Task | DP | DP3 | ACT | OpenVLA | $\pi_0$-FAST | $\pi_{0.5}$ | GR00T N1.5 |
|---|---|---|---|---|---|---|---|---|
| Articulate | Rotate chair | 100% | 90% | 100% | 70% | 100% | 100% | 100% |
| Tool Use | Use eraser to wipe the desk | 0% | 70% | 0% | 30% | 40% | 40% | 0% |
| Basic | Put dumpling toy into plate | 30% | 20% | 70% | 30% | 60% | 30% | 80% |
| Deformable | Fold towel on the desk | 0% | 20% | 0% | 40% | 20% | 40% | 50% |
| HRI | Hand over dumpling toy | 40% | 40% | 70% | 60% | 30% | 40% | 100% |
| Loco-Manip. | Walk to grab door handle | 30% | 0% | 0% | 30% | 10% | 0% | 30% |
| High Precision | Insert rose into vase | 0% | 0% | 0% | 10% | 0% | 0% | 0% |
| **Average** | | 29% | 34% | 34% | 39% | 37% | 36% | 51% |

- VLA models with pretrained priors outperform imitation learning policies.

- GR00T N1.5 achieves the best overall performance.

- All policies perform poorly on high-difficulty manipulation tasks.

# Evaluation: Cloud-based Evaluation Platform

- Website for evaluating policies trained on the *Humanoid Everyday* dataset
- Streams real robot data and records success rates
- Supports remote inference (user policy server)

# Humanoid
# Everyday

A high-frequency humanoid dataset spanning diverse everyday tasks.

● 260 Tasks   ● 7 categories   ● 30 Hz   ● 9 Modalities   ● Open Source

# Acknowledgement

Robot Learning from Any Images: Siheng Zhao, Jiageng Mao

Universal Humanoid (UH1): Jiageng Mao, Siheng Zhao

Humanoid Everyday: Hongyi Jing, Zhenyu Zhao, William Liu