# A Tale of Two Realities
## Bridging Physical Worlds with Interactable Digital Twins for Embodied Robots

### Baoxiong Jia

BIGAI

Figure generated by GPT

BIGAI

# About me
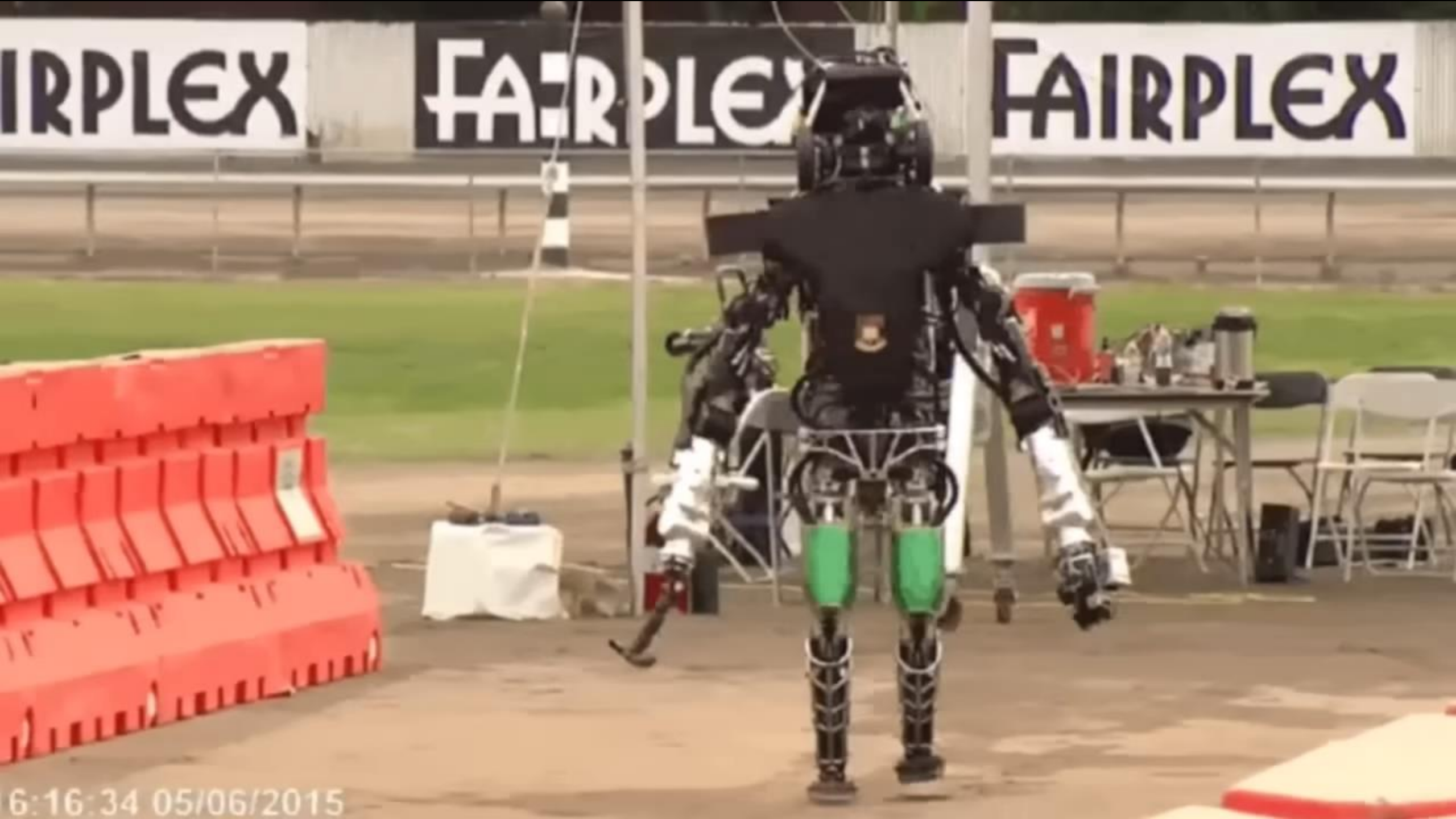
Peking University
B.S. in CS
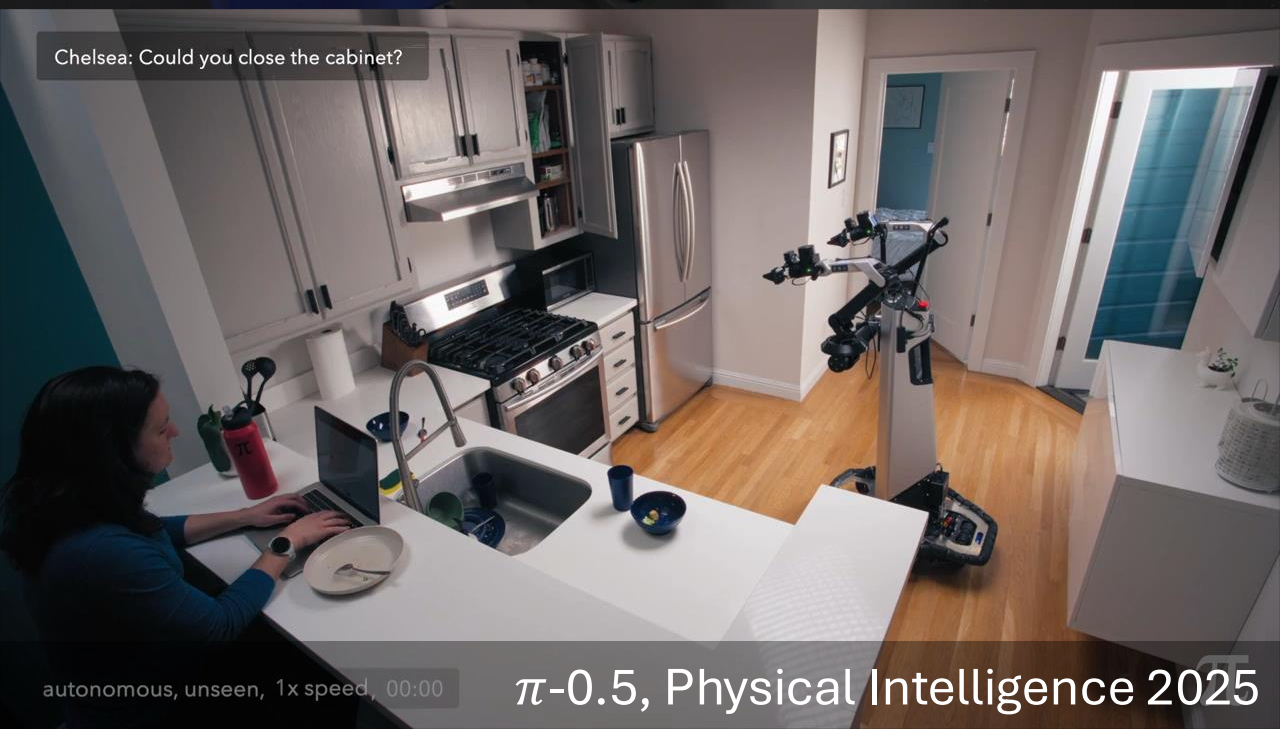2014-2018

UCLA
Ph.D. in CS
2018-2022

BIGAI
Research Scientist
2022-Present

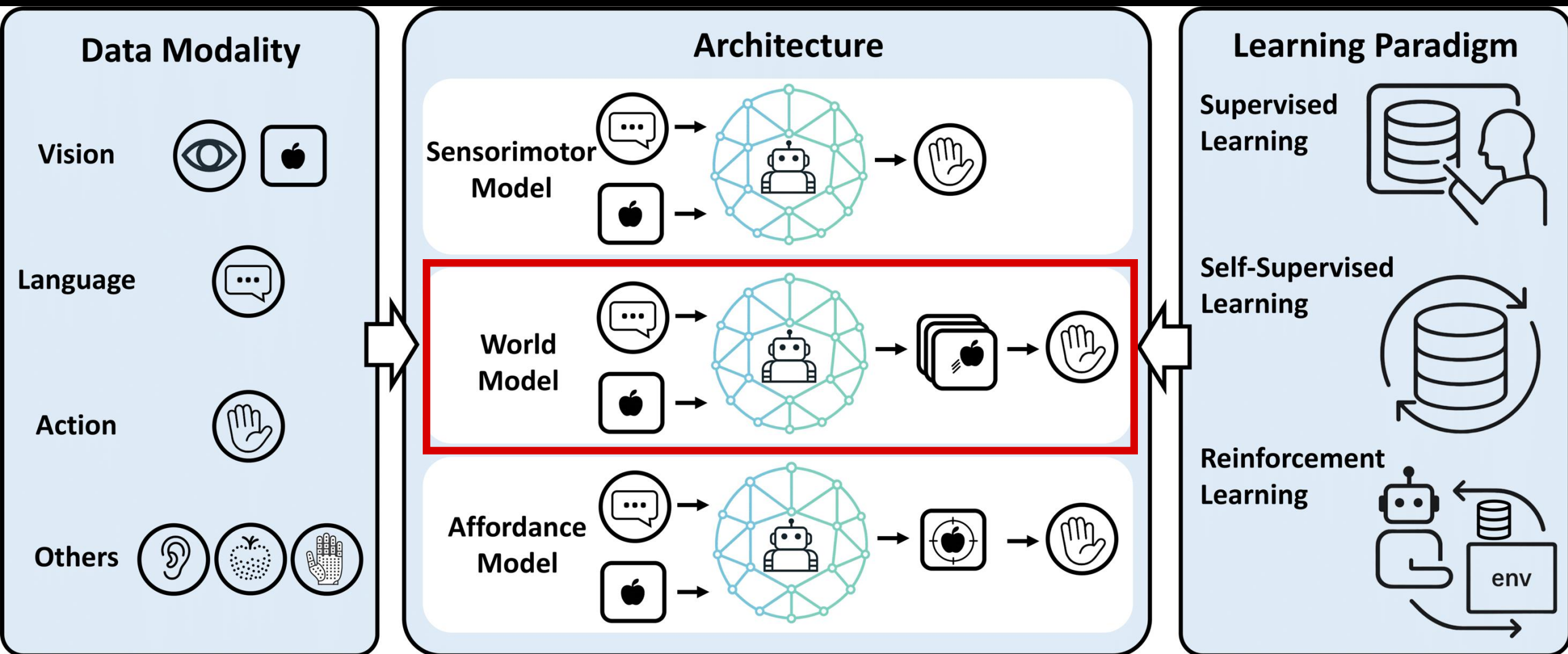Loco-Manipulation, Boston Dynamics 2025

Introducing Helix, Figure 2025

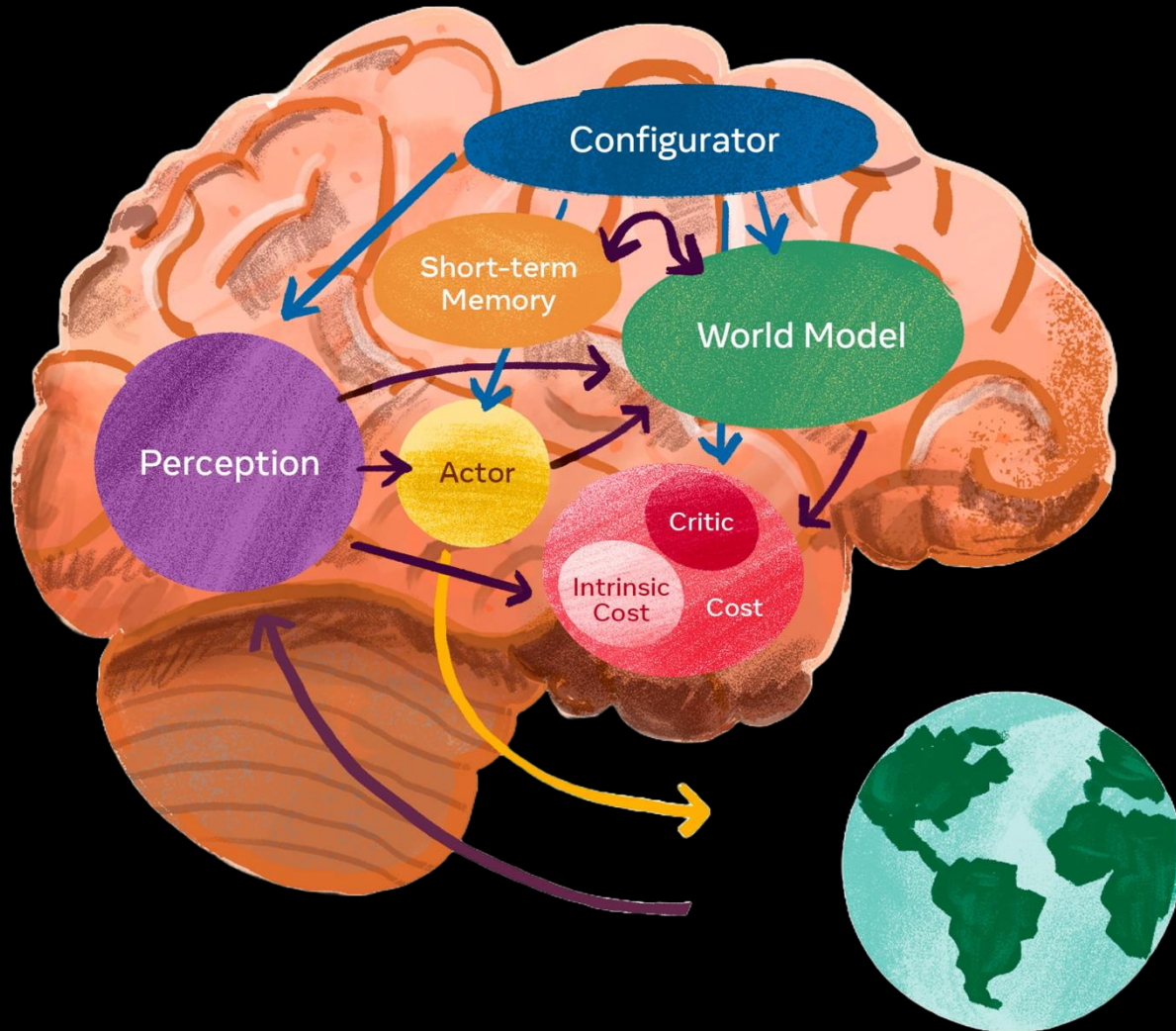Chelsea: Could you close the cabinet?

π-0.5, Physical Intelligence 2025

Early Preview of Model Capabilities, Generalist 2025

autonomous, unseen, 1x speed, 00:00

# VLA for Embodied Robots

# 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, utilise the knowledge of past events in dealing with the present and future."*

*— Kenneth Craik (1943)*

- ➢ **Integration of perception and action**
  - ❖ The model must encode states and possible actions
- ➢ **Prediction, reasoning and planning**
  - ❖ The model functions as an internal simulator for anticipating outcomes and guiding decisions
- ➢ **Efficient representation and generalization**
  - ❖ Retains essential structures to predict the future and generalize past experience
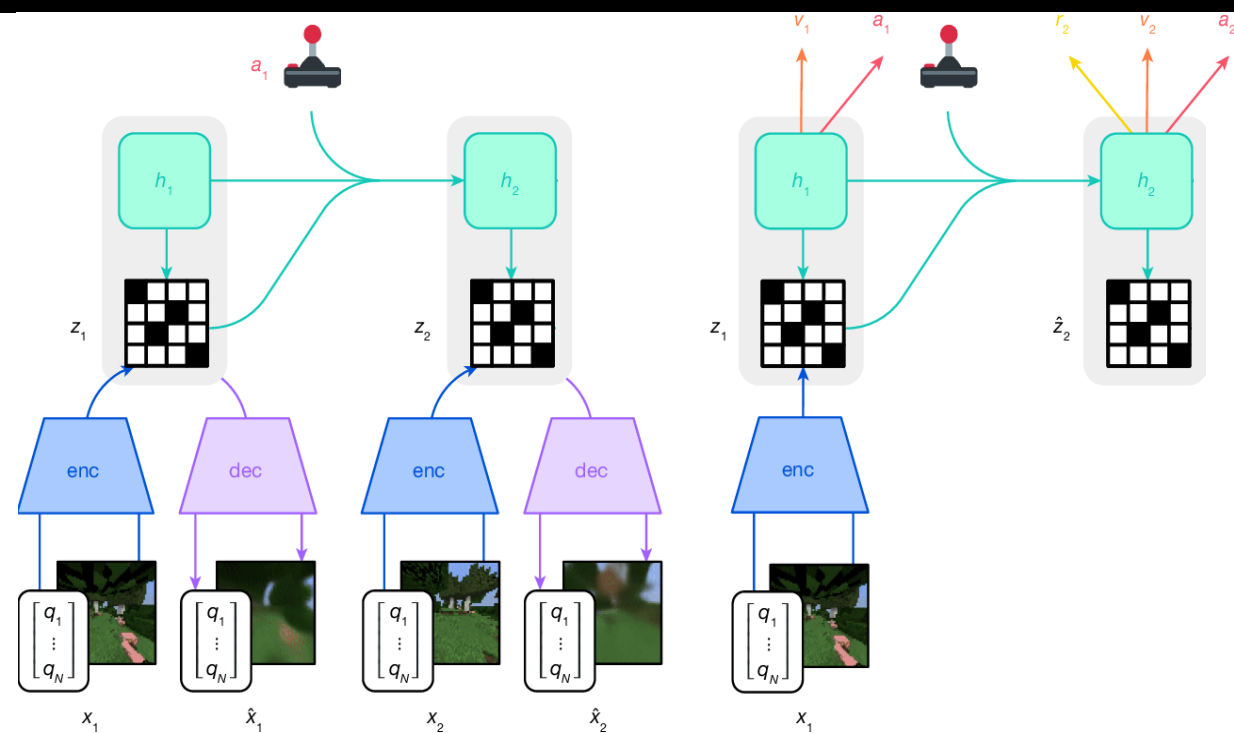
Figure credit: JEPA Blog, Meta 2022

# Model-based RL

**Representation learning for long-horizon tasks**
**Under game setting**

Dreamer 4, Google DeepMind 2025

**Model-based RL**

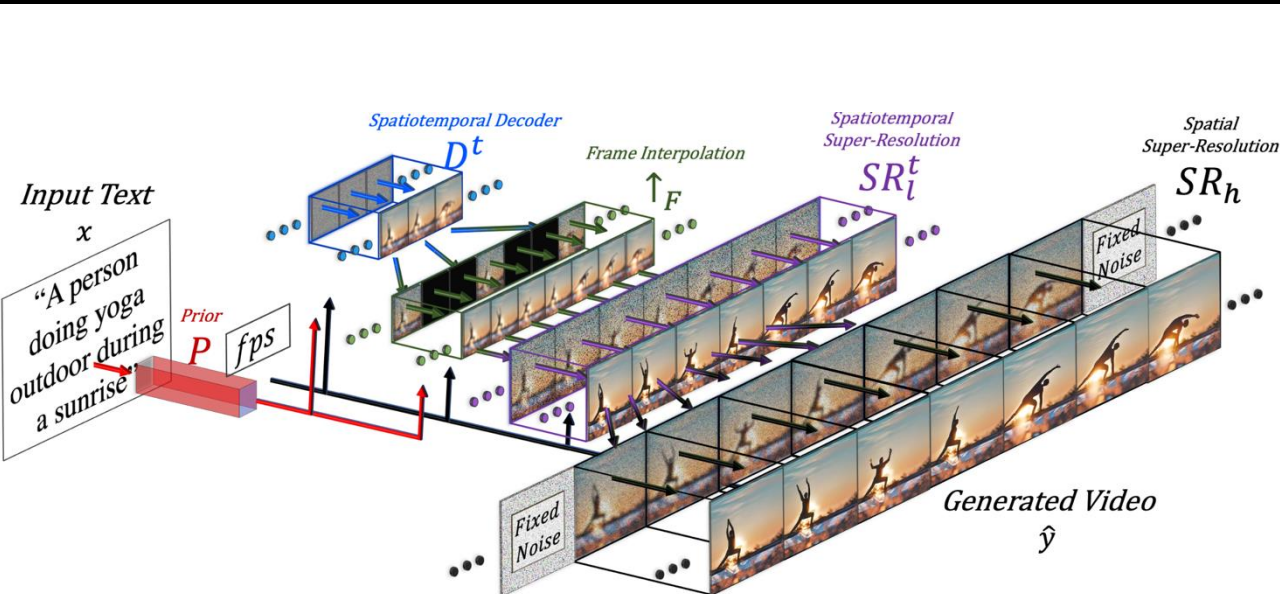Representation learning for long-horizon tasks
Under game setting

**Video Generation**

Flexible conditional generation
Weak physical consistency / modeling of action

Veo 3.1, Google Deepmind 2025

0 min

*Spatiotemporal Decoder*
$D^t$

*Frame Interpolation*
$\uparrow F$

*Spatiotemporal Super-Resolution*
$SR_l^t$

*Spatial Super-Resolution*
$SR_h$

*Input Text*
$x$

"A person doing yoga outdoor during a sunrise"

*Prior*
$P$ $fps$

*Fixed Noise*

*Generated Video*
$\hat{y}$

**Model-based RL**

Representation learning for long-horizon tasks
Under game setting

**Video Generation**

Flexible conditional generation
Weak physical consistency / modeling of action

**Latent Action Learning**

Aligning video generation with latent actions
Limited by the view point

DreamGen, NVIDIA GEAR 2025

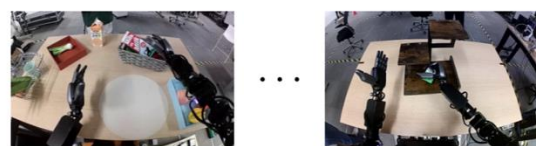Step 1. Finetune Video World Model
*Human teleoperation data*

Step 2. Rollout Video World Model
"Water the flowers"
"Pick up the tangerine"

Step 3. Label Pseudo Actions
$\hat{a}_{1:H}$   $\hat{a}_{H:2H}$
Automatically Labeled Pseudo Actions

Step 4. Visuomotor Policy Training
$\hat{a}_{1:H}$
Pseudo-labeled **neural trajectories**

**Model-based RL**

Representation learning for long-horizon tasks
Under game setting

**Video Generation**

Flexible conditional generation
Weak physical consistency / modeling of action

**Latent Action Learning**

Aligning video generation with latent actions
Limited by the view-point

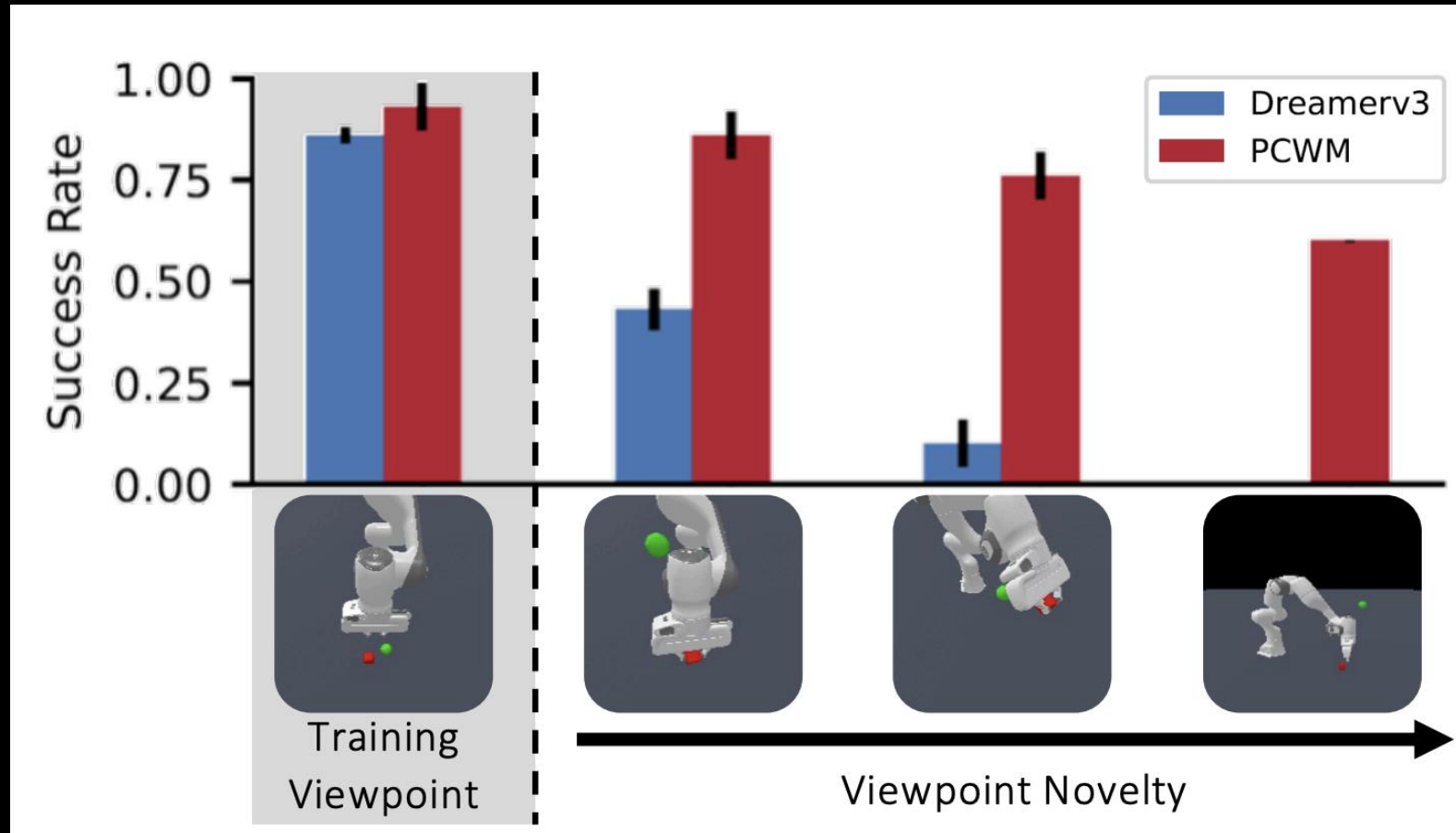**Spatial Representations**

World modeling with 3D Gaussians
Interactiveness for robot manipulation?

Marble, WorldLabs 2025

# 3D World Models?



Because of **depth estimation** challenge, tele-op must follow protocols

# 3D World Models?



Peri et al., Point Cloud Models Improve Visual Robustness in Robotic Learners, ICRA 2024

**3D helps policy learning, but requires additional sensors (RGB-D)**

# Encoding 3D Gaussians into Latent Space



(Optional)

Unposed Img.

Splatt3R

Gaussaian Splats $\mathcal{G}_t$

3D VAE

Pos. Emb.
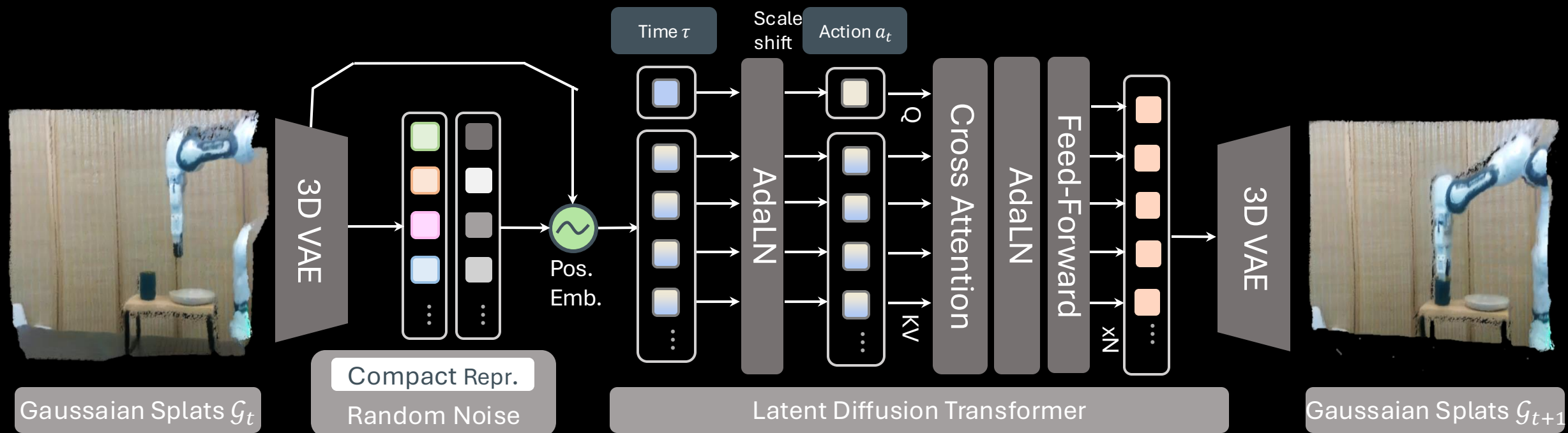
Compact Latent Representation

3D VAE

Gaussaian Splats $\mathcal{G}_t$

**Feed-Forward 3D Gaussian Reconstruction**
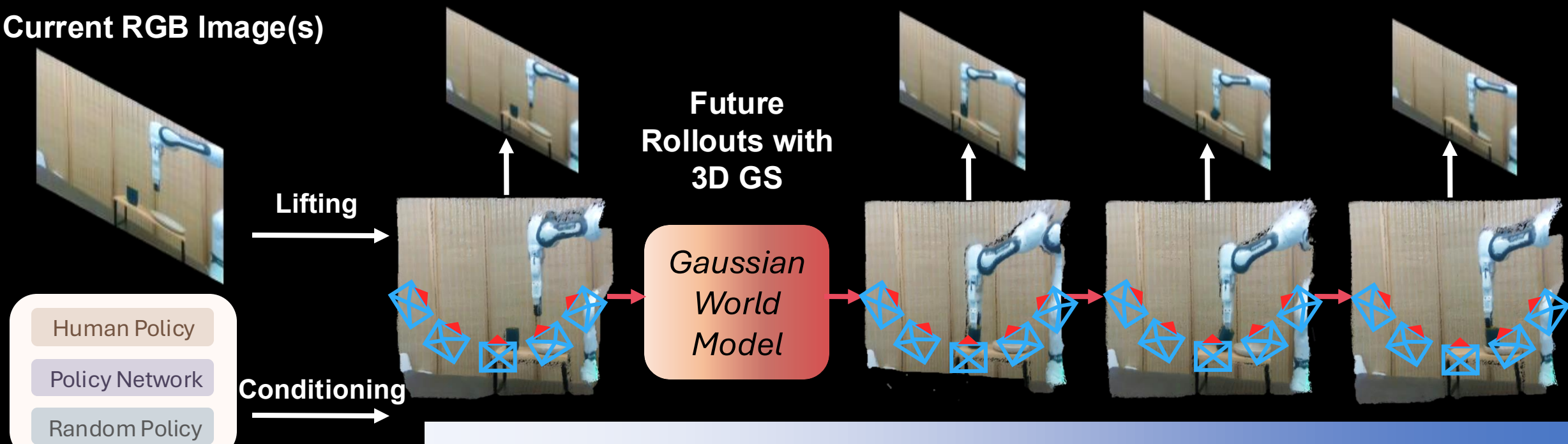
**FPS-based Subsampling Query-based Encoding**

**Rendering / Geometry Supervision**
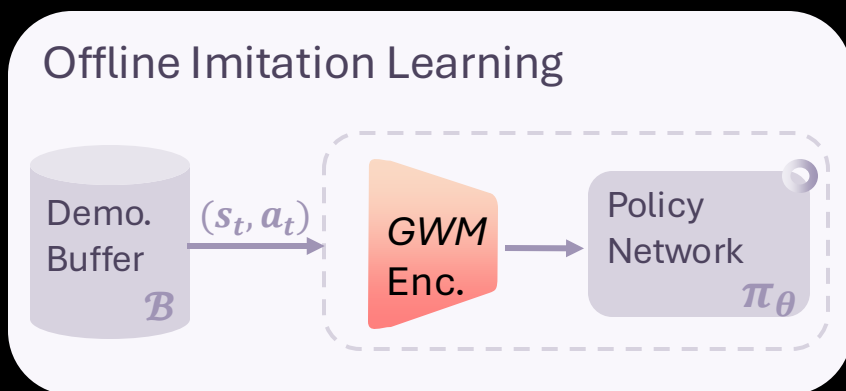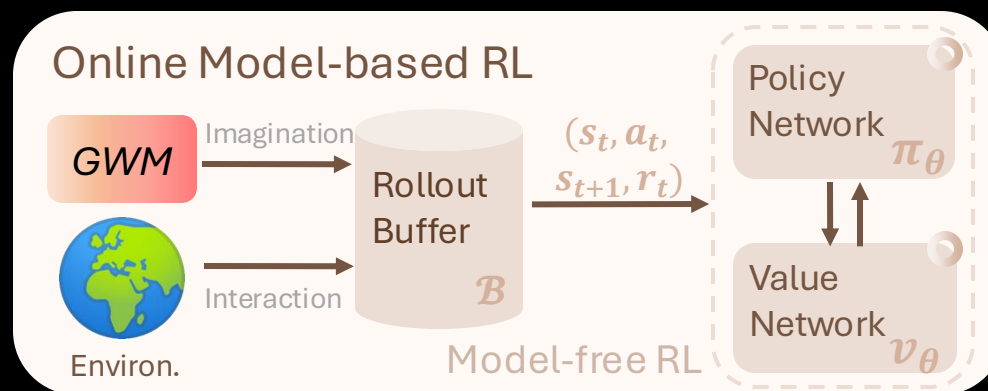
# GWM: Gaussian World Model

Gaussaian Splats $\mathcal{G}_t$

3D VAE

Compact Repr.

Random Noise

Time $\tau$

Scale shift

Action $a_t$

Pos. Emb.

AdaLN

Cross Attention

AdaLN

Feed-Forward

Q

KV

Nx

Latent Diffusion Transformer

3D VAE

Gaussaian Splats $\mathcal{G}_{t+1}$

**DiT-based Dynamics Learning and Prediction**

**Current RGB Image(s)**

Lifting

Conditioning

Human Policy

Policy Network

Random Policy

**Future Rollouts with 3D GS**

*Gaussian World Model*

Offline Imitation Learning

Demo. Buffer $\mathcal{B}$ — $(s_t, a_t)$ → GWM Enc. → Policy Network $\pi_\theta$

Online Model-based RL

*GWM* — Imagination → Rollout Buffer $\mathcal{B}$ — $(s_t, a_t, s_{t+1}, r_t)$ → Policy Network $\pi_\theta$ ⇅ Value Network $v_\theta$

Environ. — Interaction

Model-free RL

Action-conditioned (3D) Video Prediction

RGB Image → *GWM* → Future GS

Action Trajectory $\{a_t\}$

| Method | PnP CabToCounter | | PnP CounterToCab | | PnP CounterToMicrowave | | PnP CounterToSink | | PnP CounterToStove | | PnP MicrowaveToCounter | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 |
| BC-transformer | 2 | 18 | 6 | 28 | 2 | 18 | 2 | 44 | 2 | 6 | 2 | 8 |
| GWM | 18 | 32 | 4 | 22 | 14 | 44 | 20 | 38 | 2 | 18 | 20 | 26 |
| Δ | +16 | +14 | -2 | -6 | +12 | +26 | +18 | -6 | 0 | +12 | +18 | +18 |

| Method | PnP SinkToCounter | | PnP StoveToCounter | | Open SingleDoor | | Open DoubleDoor | | Close DoubleDoor | | Close SingleDoor | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 |
| BC-transformer | 8 | 42 | 6 | 28 | 46 | 50 | 28 | 48 | 28 | 46 | 56 | 94 |
| GWM | 22 | 38 | 18 | 44 | 58 | 62 | 28 | 42 | 50 | 58 | 54 | 90 |
| Δ | +14 | -4 | +12 | +16 | +12 | +12 | 0 | -6 | +22 | +12 | -2 | -4 |

| Method | Open Drawer | | Close Drawer | | TurnOn Stove | | TurnOff Stove | | TurnOn SinkFaucet | | TurnOff SinkFaucet | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 |
| BC-transformer | 42 | 74 | 80 | 96 | 32 | 46 | 4 | 24 | 38 | 34 | 50 | 72 |
| GWM | 56 | 90 | 80 | 90 | 46 | 80 | 22 | 40 | 52 | 48 | 44 | 66 |
| Δ | +14 | +16 | 0 | -6 | +14 | +24 | +18 | +16 | +14 | +14 | -6 | -6 |

| Method | Turn SinkSpout | | CoffeePress Button | | TurnOn Microwave | | TurnOff Microwave | | CoffeeServe Mug | | CoffeeSetup Mug | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 | H-50 | G-3000 |
| BC-transformer | 54 | 96 | 48 | 74 | 62 | 90 | 70 | 60 | 22 | 34 | 0 | 12 |
| GWM | 72 | 90 | 76 | 90 | 64 | 84 | 70 | 54 | 36 | 50 | 16 | 28 |
| Δ | +18 | -6 | +28 | +16 | +2 | -6 | 0 | -6 | +14 | +16 | +16 | +16 |

# GWM for Online Model-based RL



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

# GWM for Real-World Robot Manipulation



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

# Takeaways

➢ **Encoding explicit spatial information into world modeling**

  ❖ **Unite world modeling with 3D generation, video generation, multi-view reconstruction, etc.**

# Takeaways

➢ Encoding explicit spatial information into world modeling

❖ Unite world modeling with 3D generation, video generation, multi-view reconstruction, etc.

➢ **Better VLA modeling with world modeling**

❖ **Using the latent representation alone does not fully utilize the predictive power of world models**

# Takeaways

➢ Encoding explicit spatial information into world modeling

  ❖ Unite world modeling with 3D generation, video generation, multi-view reconstruction, etc.

➢ Better VLA modeling with world modeling

  ❖ Using the latent representation alone does not fully utilize the predictive power of world models

➢ **Scalable 4D world modeling**

  ❖ **Scalability vs. precision still stands as an issue, feed-forward 3D Gaussians still need improvement**

  **Especially for Dynamic Objects**

# Manipulation Policies involve Dynamic Objects



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

# Provide motion prior from pre-trained tracking models

**Key Insights:** Analyze noisy 3D tracks to provide robust initialization and optimization signals



Filtering noise and estimate articulation parameters

# Camera, Depth, Tracks Estimation

Video Frames



$I_0$          $I_t$

# Camera, Depth, Tracks Estimation

# Camera, Depth, Tracks Estimation

# Motion Analysis & Deformation Field Initialization

# Motion Analysis & Deformation Field Initialization

# Motion Analysis & Deformation Field Initialization

# Motion Analysis & Deformation Field Initialization

# Motion Analysis & Deformation Field Initialization

# Geometry Reconstruction & Articulation Learning

# Geometry Reconstruction & Articulation Learning

# Geometry Reconstruction & Articulation Learning

Video Frames

$I_0$    $I_t$
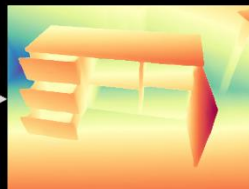
VGGT

Depth & Pose

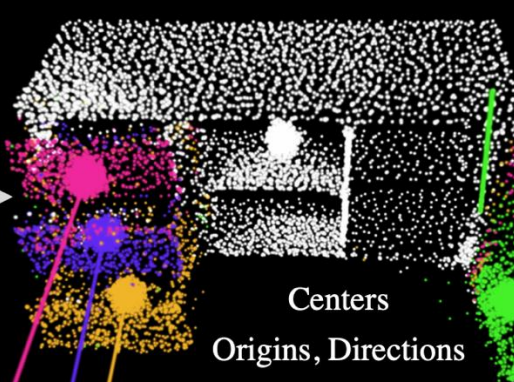$D_0$    $D_t$

TAPIP3D

3D Tracks

$x^0$   Sample   $x^t$

3D Tracks → Motion Analysis → Static, Noise, Prismatic, Revolute → Motion Clustering → Centers Origins, Directions → Init → Deformation Field

$x^{t_1}$   $x^{t_2}$

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

$\hat{x}^{t_2} = \mathcal{F}(\hat{x}^c, t_2)$
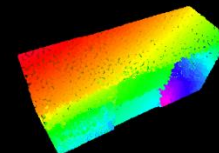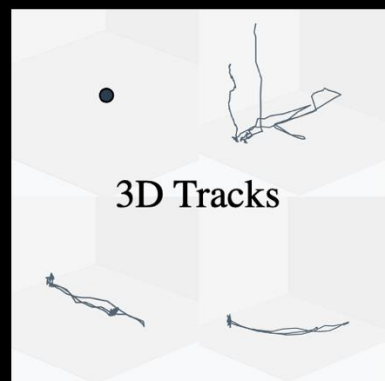
First N-frames → Init → Canonical Gaussians $\mathcal{G}^c$ → Deformation Field (Articulation-based Deformation) → $\mathcal{G}^t$ → Deformed Gaussians → Render

$\hat{I}_t \longleftrightarrow I_t$

$\hat{D}_t \longleftrightarrow D_t$

# Quantitative Comparison

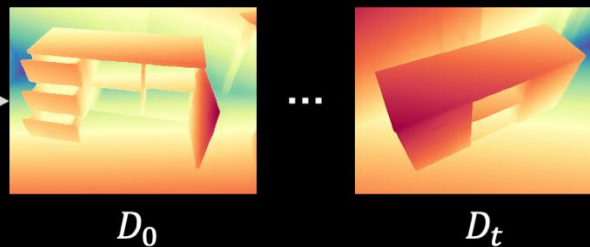| Method | Revolute Joint Estimation | | | Prismatic Joint Estimation | | Reconstruction | | |
|---|---|---|---|---|---|---|---|---|
| | Axis (○) | Position (cm) | State (○) | Axis (deg) | State (cm) | CD-w (cm) | CD-m (cm) | CD-s (cm) |
| ArticulateAnything[†] (Le et al., 2025) | 46.98±45.27 | 81.00±40.00 | N/A | 52.71±44.69 | N/A | 11.00±22.00 | 59.00±73.00 | 7.00±18.00 |
| RSRD[†] (Kerr et al., 2024) | 67.06±29.22 | 203.00±748.00 | 59.02±34.38 | 69.91±24.07 | 70.00±48.00 | 339.00±2147.00 | 82.00±117.00 | 14.00±41.00 |
| Video2Articulation[†] (Peng et al., 2025) | 18.34±32.09 | 13.00±25.00 | 14.32±26.35 | 13.75±18.91 | 8.00±22.00 | 1.00±1.00 | 13.00±26.00 | 6.00±19.00 |
| Video2Articulation (Peng et al., 2025) | 13.83±28.15 | 11.55±22.39 | 10.25±21.27 | 14.37±19.08 | 3.44±6.25 | 3.45±16.46 | 12.21±24.44 | 5.39±17.09 |
| **Ours** | **0.32±0.44** | **0.42±0.75** | **1.15±2.29** | **0.35±0.45** | **1.03±2.46** | **0.29±0.24** | **0.40±0.32** | **1.11±2.11** |

| Method | Axis (°) | Position(cm) | CD-w(cm) | CD-m(cm) | CD-s(cm) |
|---|---|---|---|---|---|
| ArticulateAnything (Le et al., 2025) | 43.65 ± 44.72 | 15.66 ± 36.20 | 16.10 ± 37.34 | 17.66 ± 36.74 | 16.04 ± 37.36 |
| Video2Articulation (Peng et al., 2025) | 48.88 ± 24.18 | 37.04 ± 31.82 | 5.07 ± 21.78 | 30.63 ± 25.64 | 10.22 ± 22.23 |
| **Ours** | **0.34±0.80** | **0.10±0.10** | **0.09±0.09** | **0.26±0.61** | **0.24±0.58** |

**State-of-the-art performance on all metrics**

**Reducing the error by about two orders of magnitude**

# Qualitative Comparison



| Input Frames | Articulate Anything | Video2Articulation | Ours | GT |

# Real-world Experiments

Data Capture

Input

Recon

Gaussian

# Takeaways

➢ **Utilizing motion priors is crucial for dynamic object modeling**

❖ **Articulated objects are still easy to model, priors or patterns are more difficult to be defined**

# Takeaways

➢ Utilizing motion priors is crucial for dynamic object modeling
   ❖ Articulated objects are still easy to model, priors or patterns are more difficult to be defined

➢ **Object-level articulated object reconstruction is do-able**
   ➢ Generating an interactable scene is still very difficult, due to both the increasing number of dynamic parts and occlusions
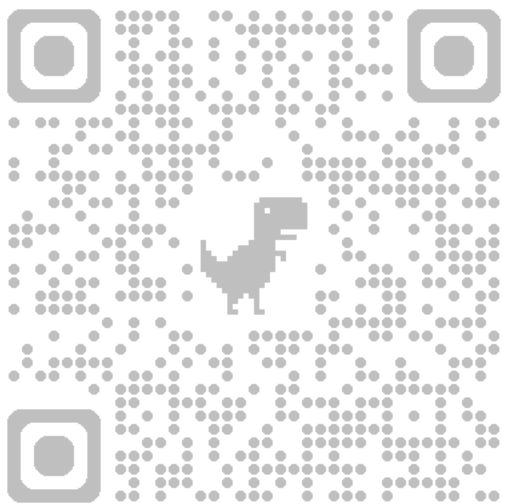
# Takeaways

➤ **Utilizing motion priors is crucial for dynamic object modeling**

  ❖ **Articulated objects are still easy to model, priors or patterns are more difficult to be defined**

➤ **Object-level articulated object reconstruction is do-able**

  ❖ **Generating an interactable scene is still very difficult, due to both the increasing number of dynamic parts and occlusions**

➤ **Monocular video with sufficient camera trajectory design gives good reconstruction results**

  ❖ **How to utilize large-scale internet-scale egocentric interaction data remains a challenge**

GWM: Towards Scalable Gaussian World
Models for Robotic Manipulation
ICCV 2025

https://gaussian-world-model.github.io/

**Thank you**
**Q&A**

COLA (arXiv 2025)

SceneWeaver (NeurIPS 2025)

VideoArtGS: Building Digital Twins of
Articulated Objects from Monocular Video
arXiv:2509.17647
https://videoartgs.github.io