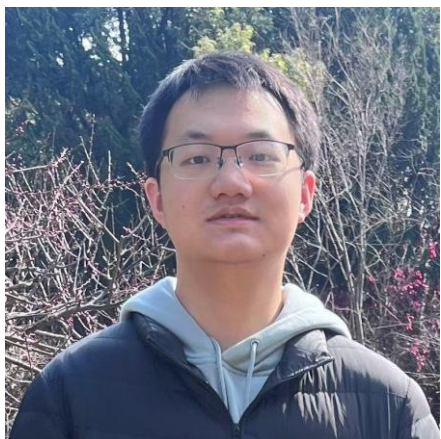


# VisualMimic:

## Visual Humanoid Loco-Manipulation via Motion Tracking and Generation



Shaofeng Yin\*



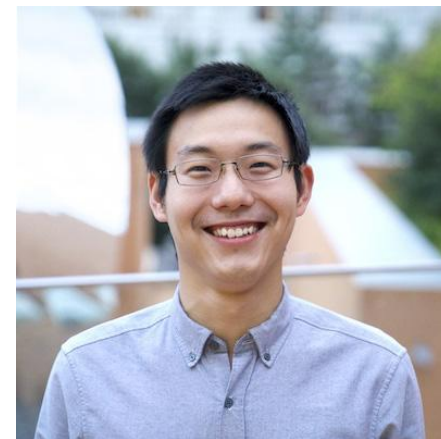
Yanjie Ze\*



Koven Yu



C. Karen Liu†



Jiajun Wu†

\*Contributed Equally    † Advised Equally



# Motivation: Visual Humanoid Control



**Project Goal:** Develop a **general\*** **sim-to-real** **visual** **whole-body control** framework for humanoid loco-manipulation.

\* general: capable of performing diverse tasks; easy to train and add new tasks

# Motivation: Visual Humanoid Control



**Project Goal:** Develop a **general\*** **sim-to-real** **visual** **whole-body control** framework for humanoid loco-manipulation.

\* general: capable of performing diverse tasks; easy to train and add new tasks



# Motivation: Visual Humanoid Control



**Project Goal:** Develop a **general\*** **sim-to-real** **visual** **whole-body control** framework for humanoid loco-manipulation.

\* general: capable of performing diverse tasks; easy to train and add new tasks

# Motivation: Visual Humanoid Control



**Project Goal:** Develop a **general\*** **sim-to-real** **visual** **whole-body control** framework for humanoid loco-manipulation.

\* general: capable of performing diverse tasks; easy to train and add new tasks

# Motivation: **Visual** Humanoid Control

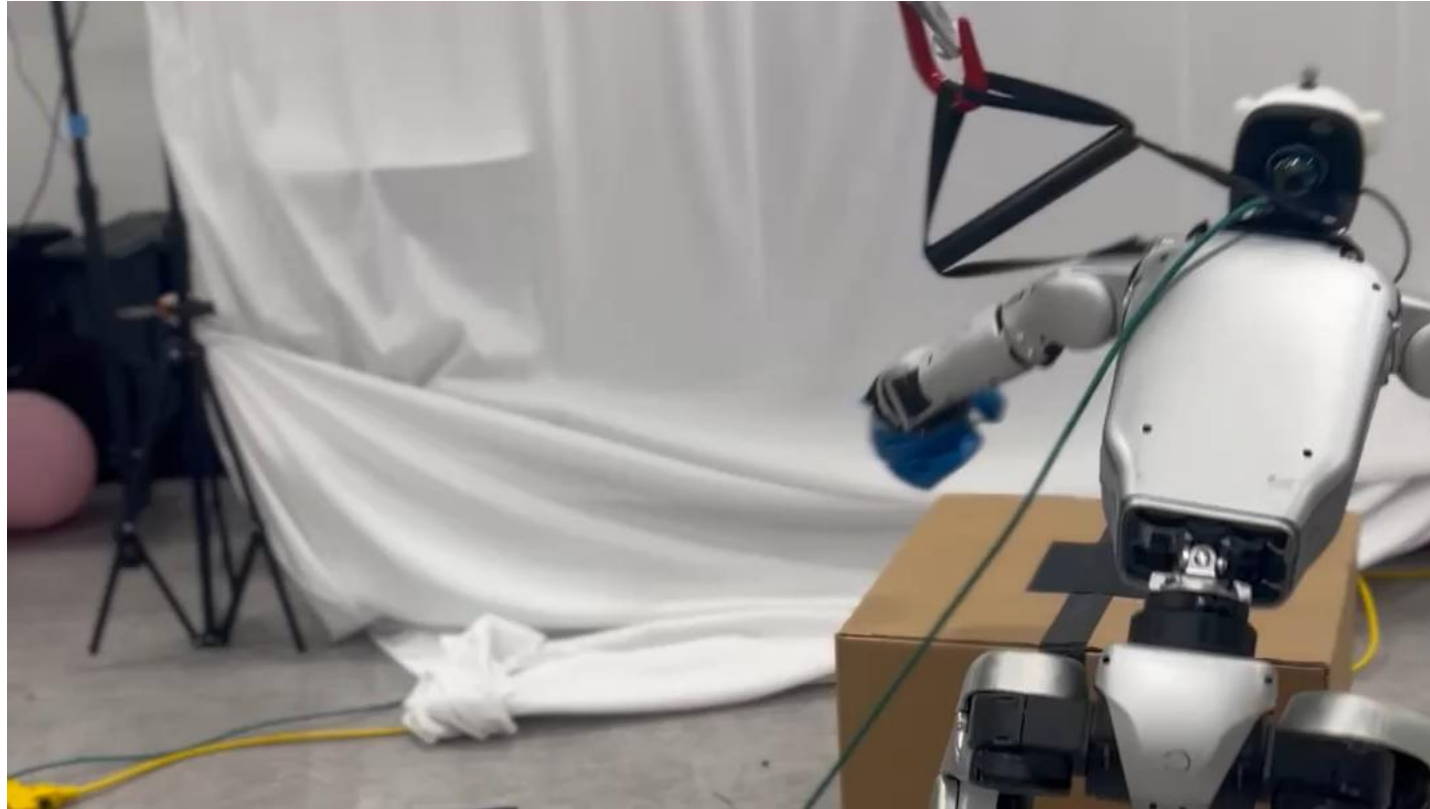


**Project Goal:** Develop a **general\*** **sim-to-real** **visual** **whole-body control** framework for humanoid loco-manipulation.

\* general: capable of performing diverse tasks; easy to train and add new tasks



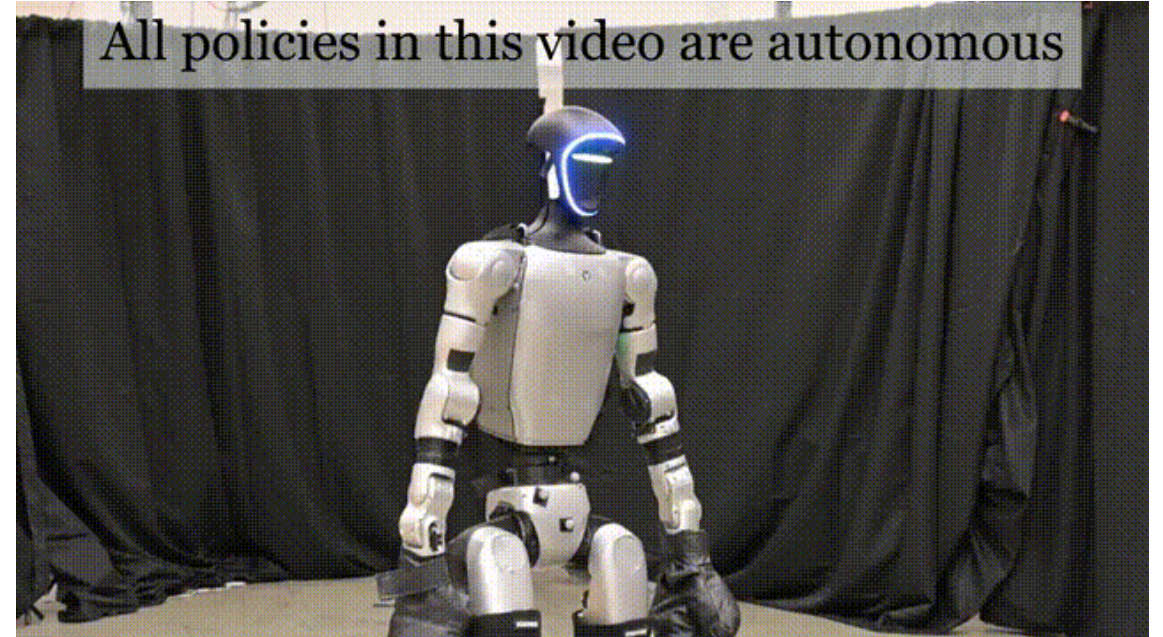
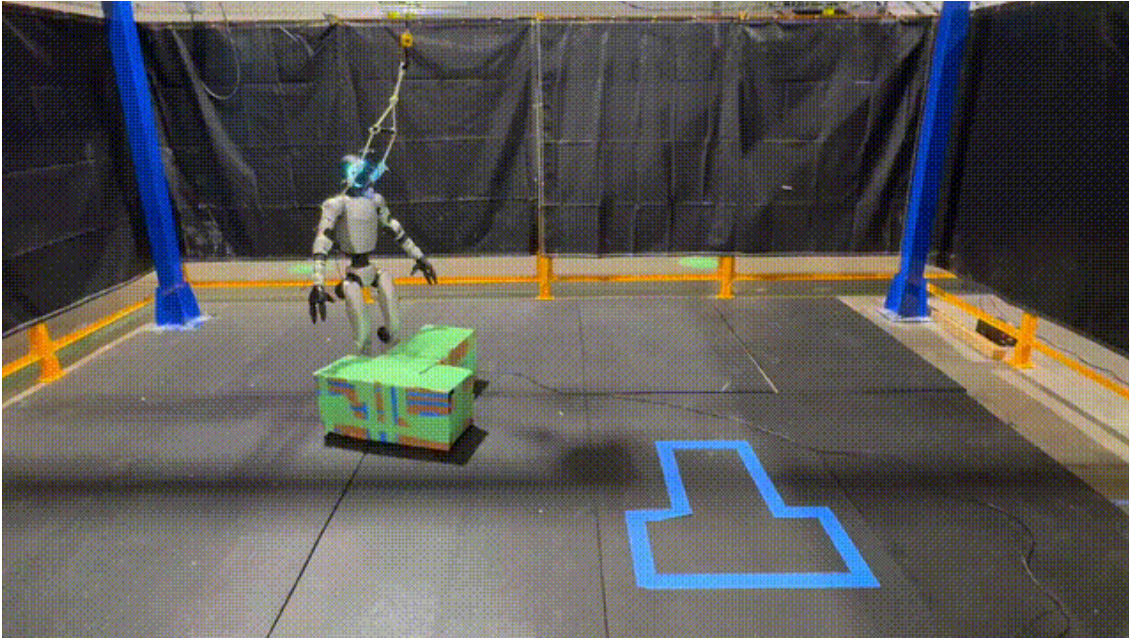
# Motivation: Robust Visual Humanoid Control



**Project Goal:** Develop a **general\*** **sim-to-real** **visual** **whole-body control** framework for humanoid loco-manipulation.

\* general: capable of performing diverse tasks; easy to train and add new tasks

# Motivation: **Robust** Visual Humanoid Control



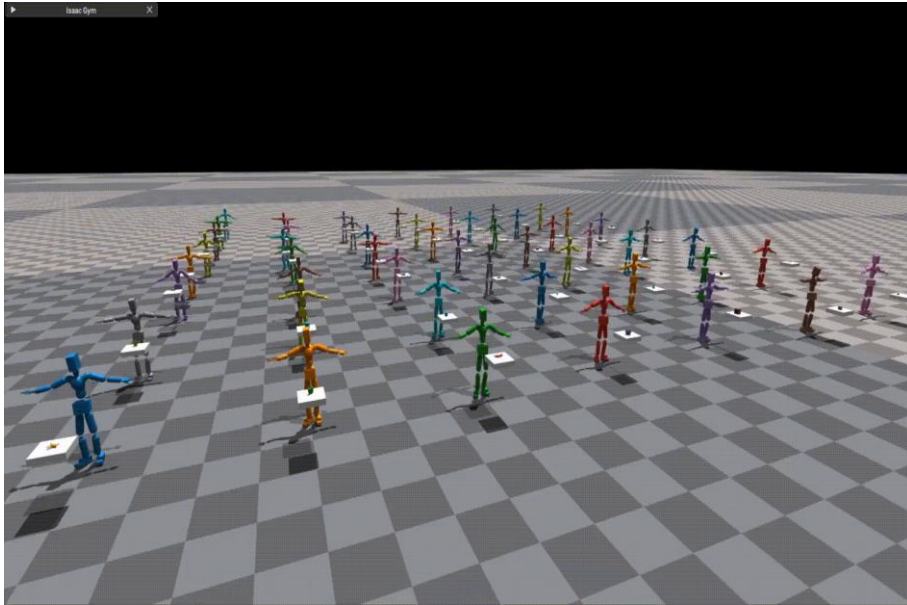
## Other Approaches

Imitation Learning

“Imitation Learning”

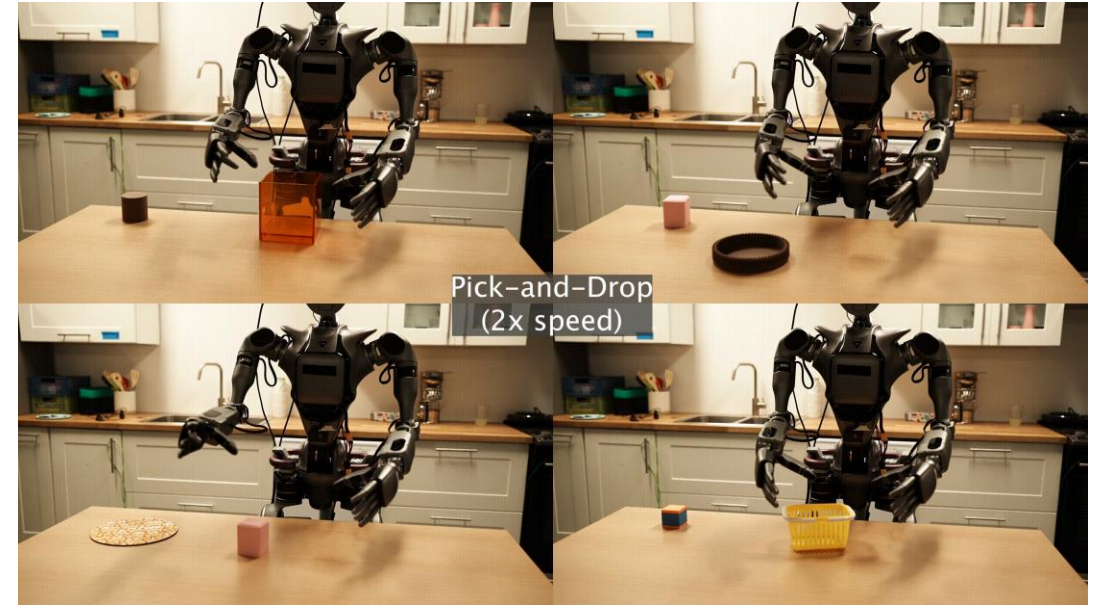


# Two Paradims in Humanoid Policy Training



train end-to-end policy  
with task-specific human motions

✗ need task-relevant motion



trained end-to-end policy  
with complex RL reward design

✗ complex reward design

## Our Project

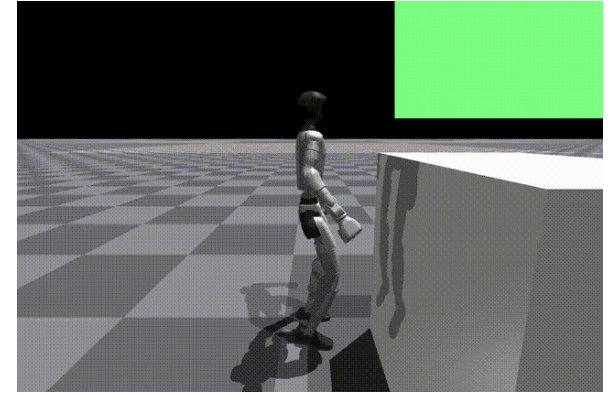
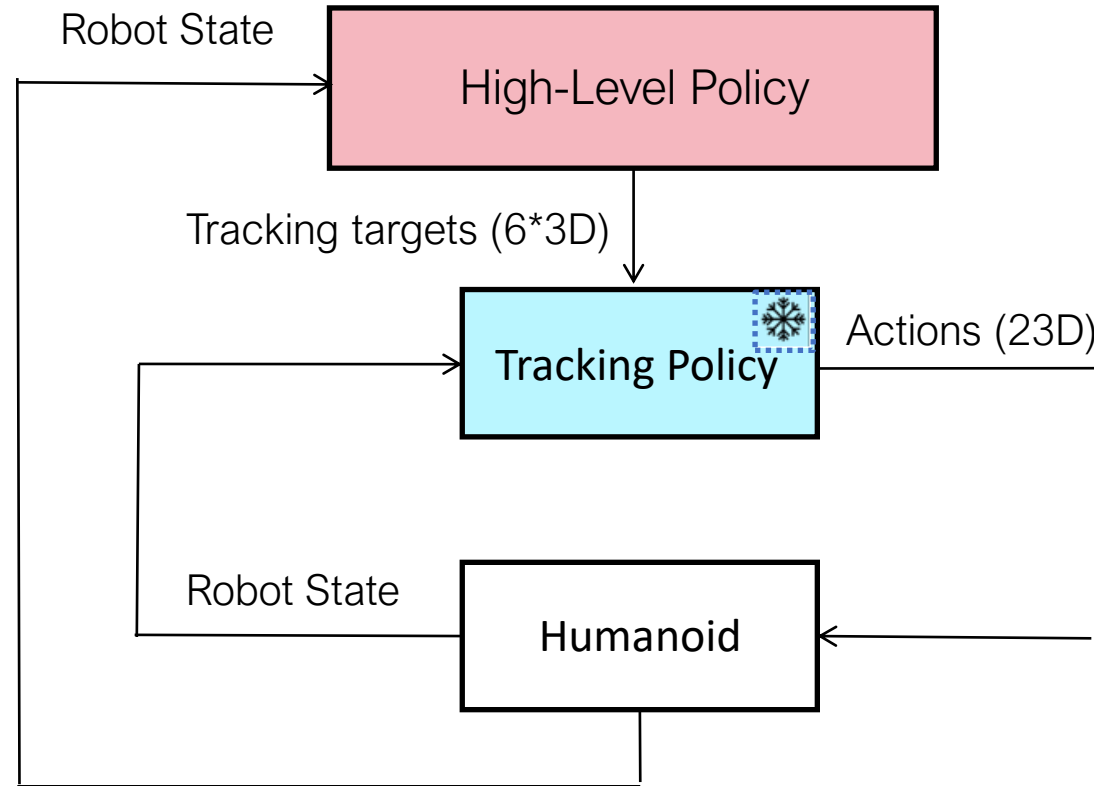
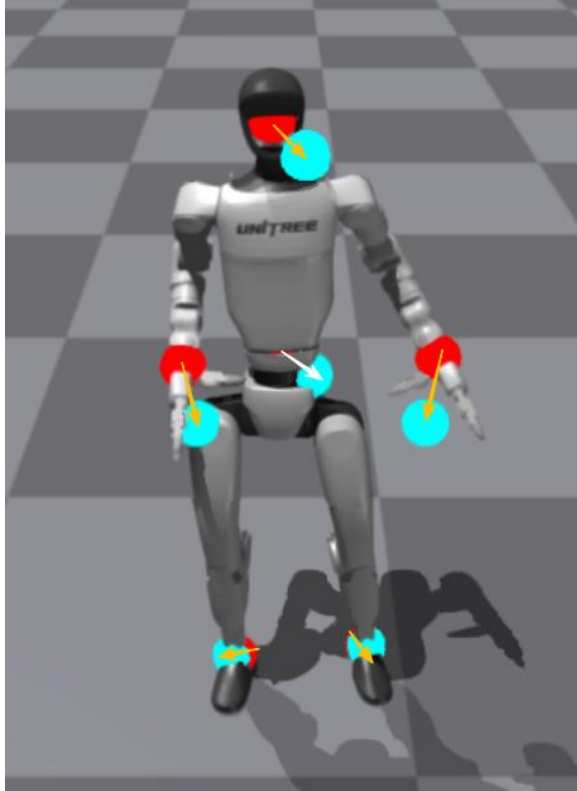
train a low-level policy with task-agnostic human motions  
train high-level policies with simple RL reward

- ✓ no need for task-relevant motion
- ✓ easy reward design
- ✓ quickly adapt to a new task

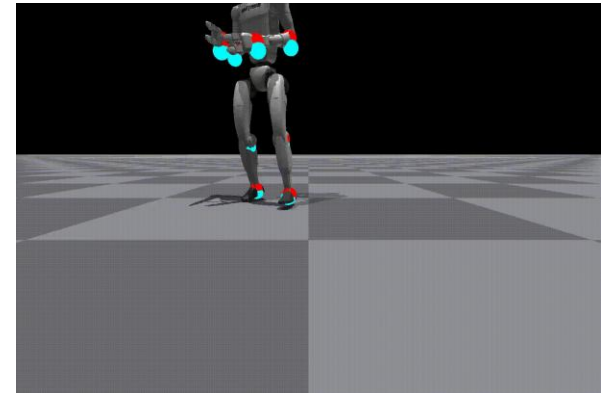
Omnigrasp: Grasping Diverse Objects with Simulated Humanoids, Luo et al., 2024

Sim-to-Real Reinforcement Learning for Vision-Based Dexterous Manipulation on Humanoids, Lin et al., 2025

# Method: Hierarchical Framework



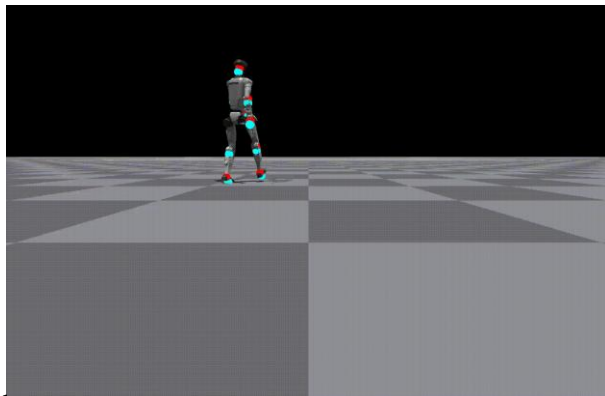
High-Level: Keypoint Generator



Low-level: General Keypoint Tracker

A hierarchical framework with a flexible interface that enables the high-level policy to control all body parts

# Low-Level: General Keypoint Tracker (Stage 0)

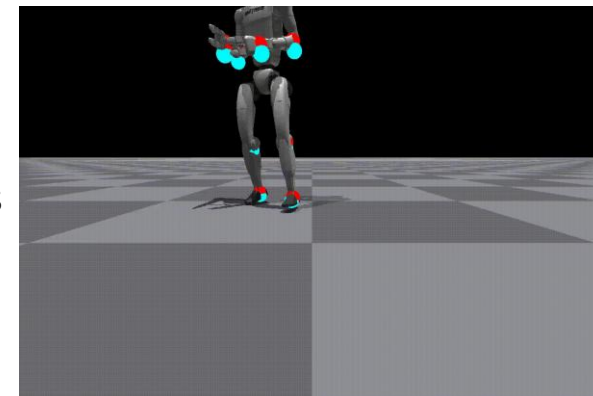


distill

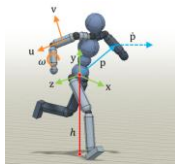
whole-body  
motion  
commands



6 keypoints  
commands



20 future steps of target motion



proprioceptive observation

General Keypoint Tracker

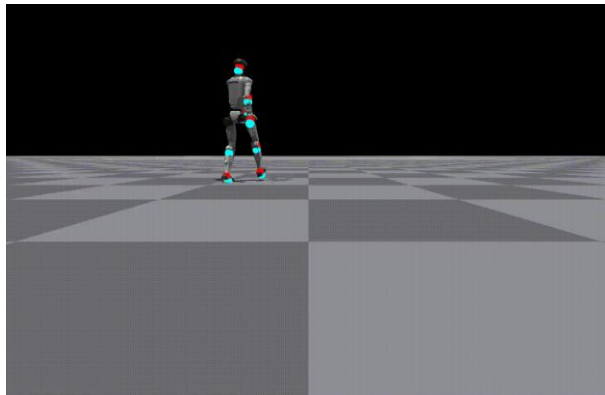
PD controller commands

Learned through RL using motion tracking rewards

(Stage 0)



# Low-Level: General Keypoint Tracker (Stage 1)

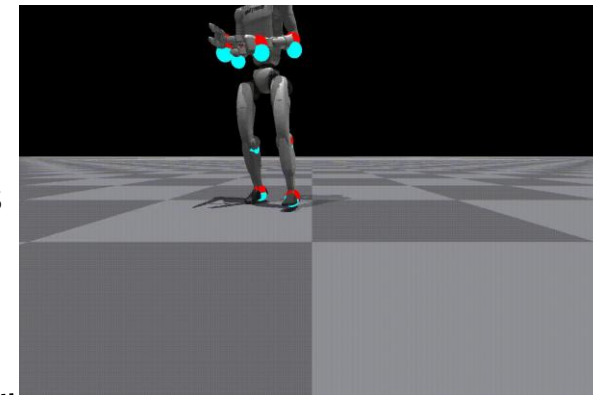


distill

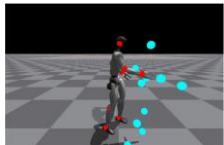
whole-body  
motion  
commands



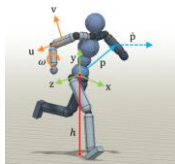
6 keypoints  
commands



(Stage 1)



**one** future step of target motion  
(simplified command)



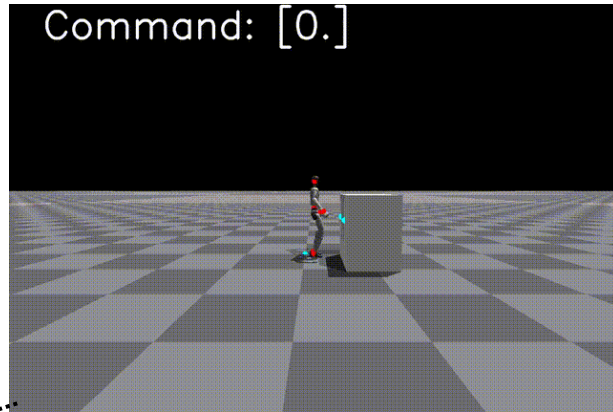
proprioceptive observation (**partial**)


General Keypoint Tracker

PD controller commands

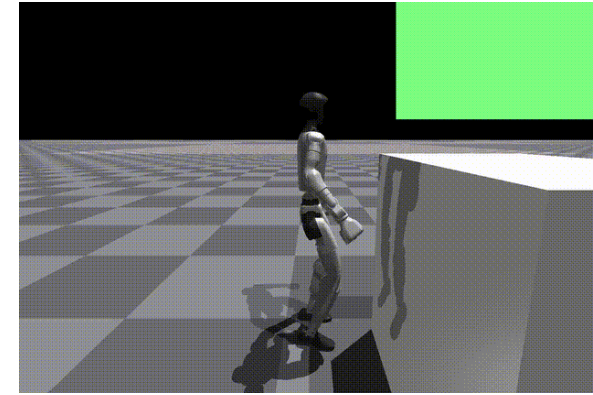
Learned through **BC** with **D**Agger

# High-Level: Task-Specific Keypoint Generator (Stage 2)

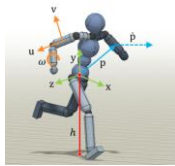


state input  visual input

distill



state-based task-relevant input  
e.g. target point position  
(10.0, 0.0, 0.8)



proprioceptive observation

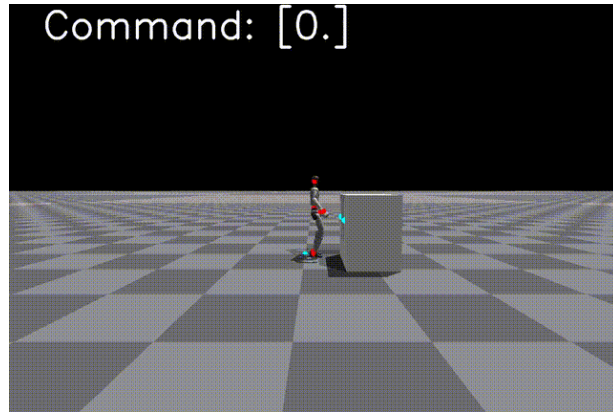
Task-Specific Keypoint  
Generator

Tracking target for low-level

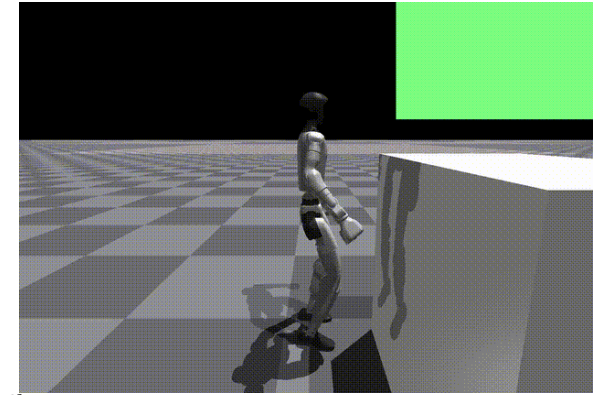
Learned through RL using task-specific rewards

(Stage 2)

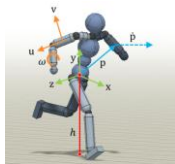
# High-Level: Task-Specific Keypoint Generator (Stage 3)



state input  $\xrightarrow{\text{distill}}$  visual input



Depth image



proprioceptive observation (partial)

Task-Specific Keypoint Generator

Tracking target for low-level

(Stage 3)

Learn through BC with DAgger



# High-Level: Reward Design

- Approach

- $R_{\text{approach}}(t) = e^{-0.1d(t)}$

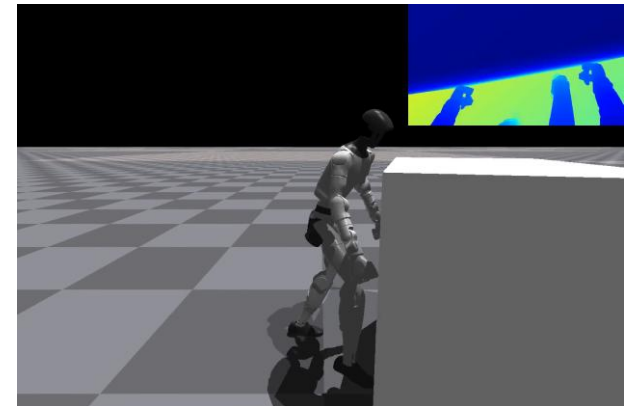
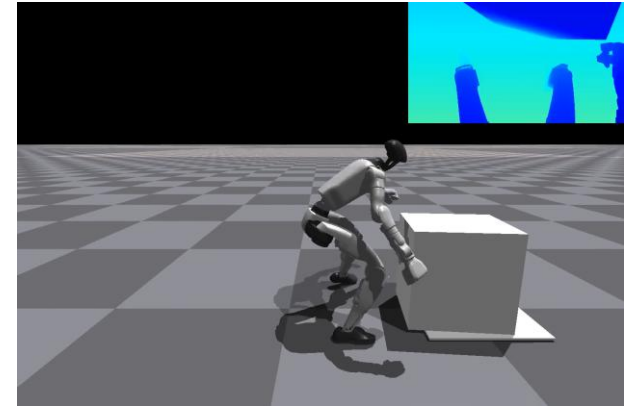
- $R_{\text{approach}}(t) = \frac{2 e^{-0.1d_1(t)} e^{-0.1d_2(t)}}{e^{-0.1d_1(t)} + e^{-0.1d_2(t)}}$

- Forward

- $R_{\text{forward}}(t) = \tanh\left(10[x_{\text{obj}}(t) - \max_{t' < t} x_{\text{obj}}(t')]\right)_+$

- Force

- $R_{\text{force}}(t) = e^{-0.1[F_{\text{des}} - F_{\text{obj}}(t)]_+}$



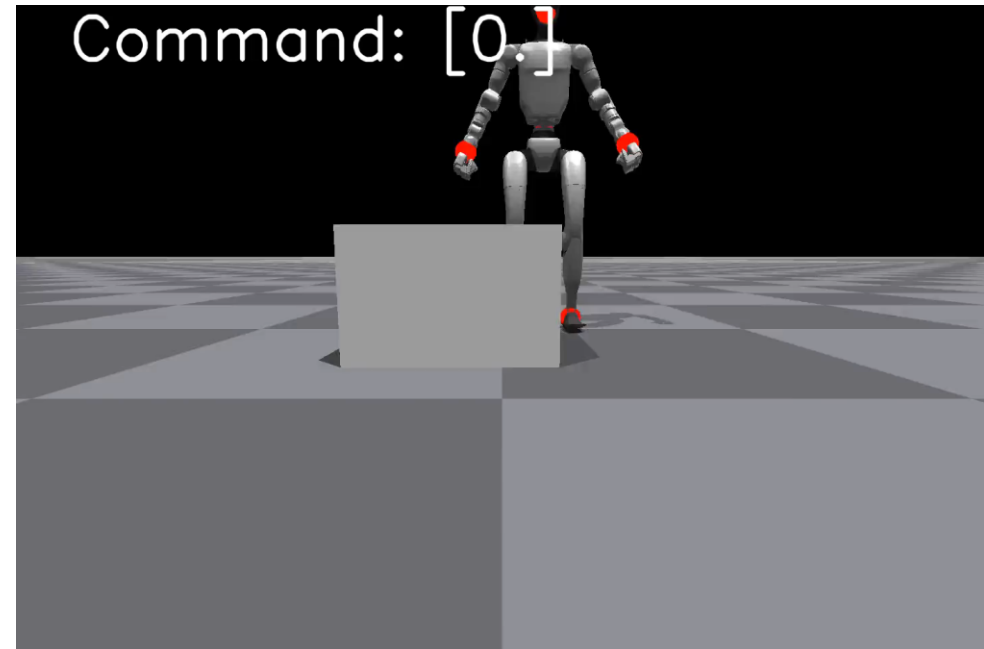
# High-Level: Reward Design

- Look at object

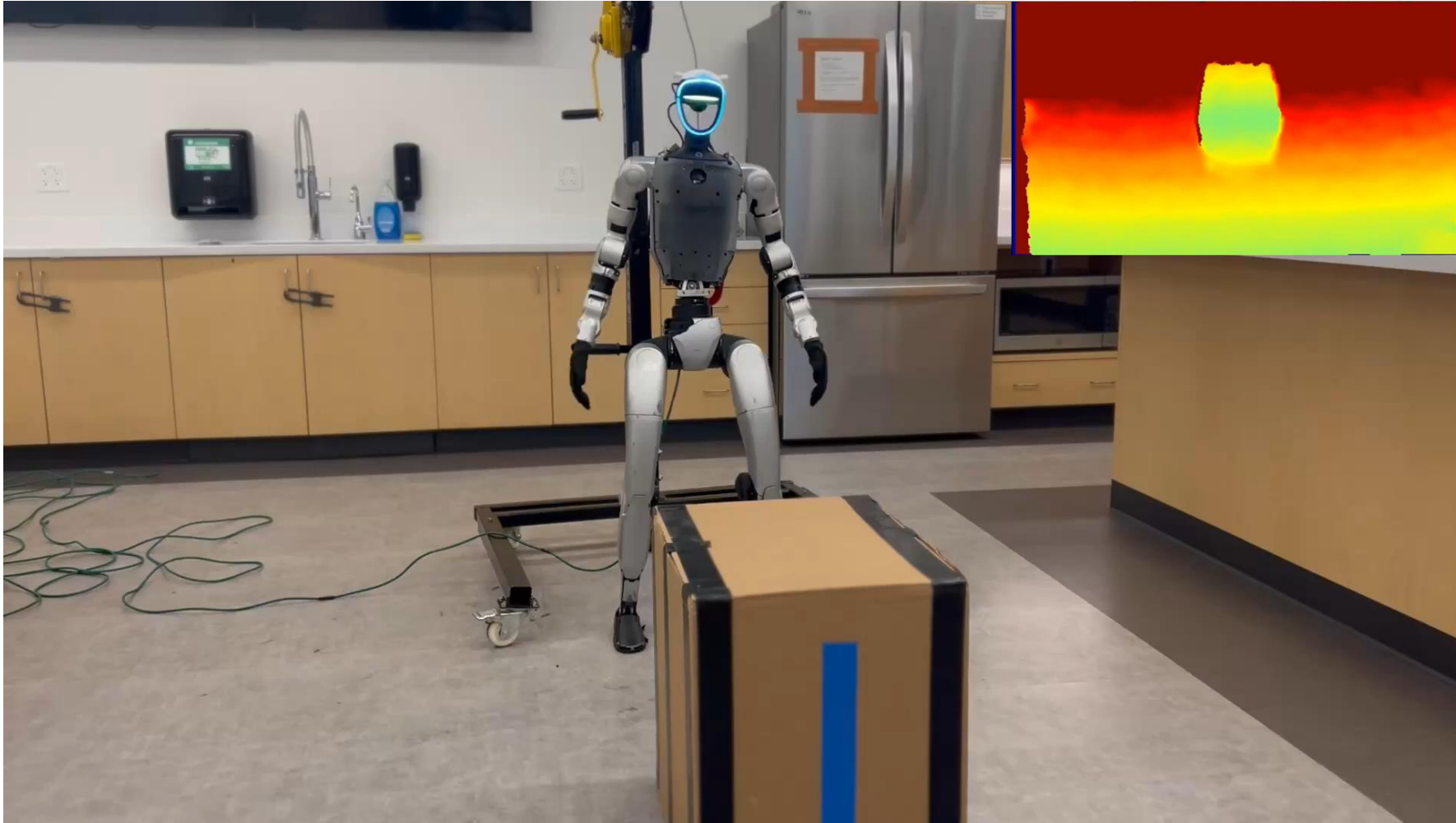
- $R_{\text{look}}(t) = -(\arccos(\hat{\mathbf{f}}_{\text{body}} \cdot \hat{\mathbf{d}}_{\text{obj}}))^2$

- Drift

- $R_{\text{drift}}(t) = -\tanh\left(10[|y_{\text{obj}}(t)| - \max_{t' < t} |y_{\text{obj}}(t')|]_+\right).$



# Real World: Visual Mask



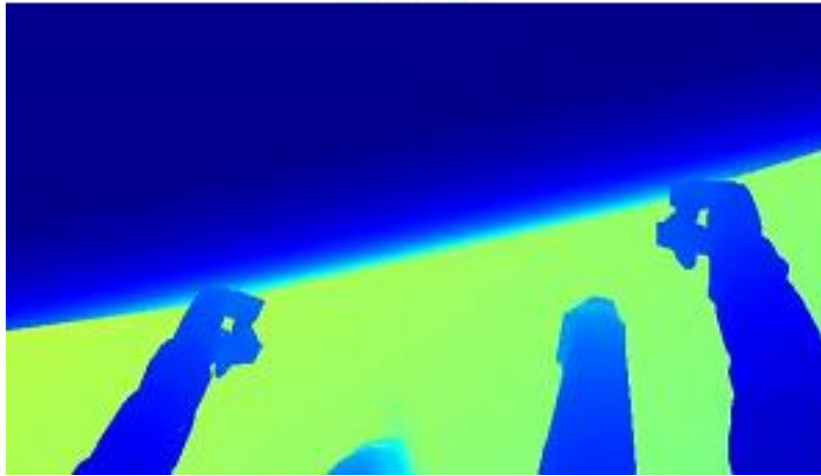


# Real World: Visual Mask

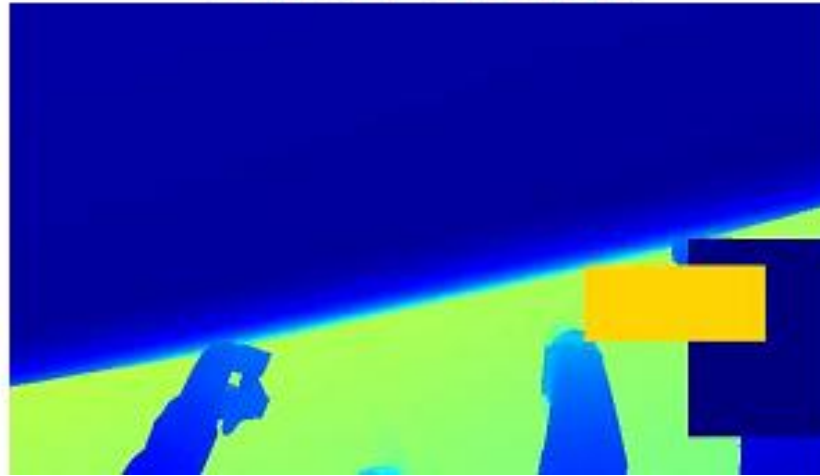


# Real World: Visual Mask

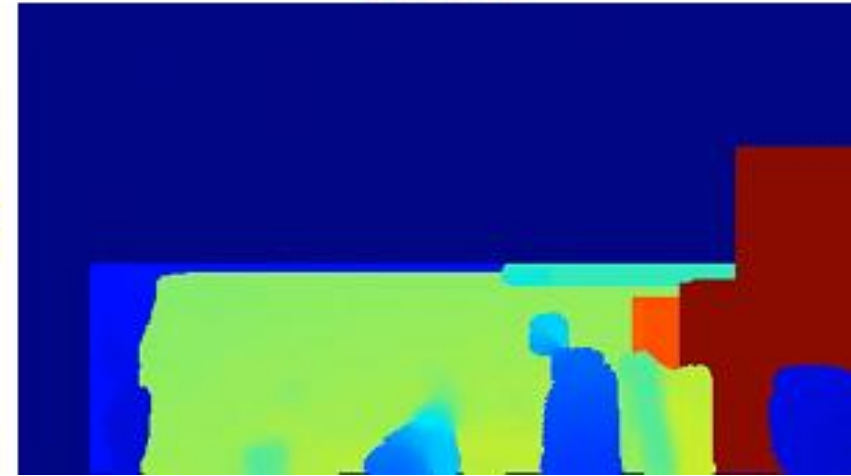
Sim



Sim + Mask

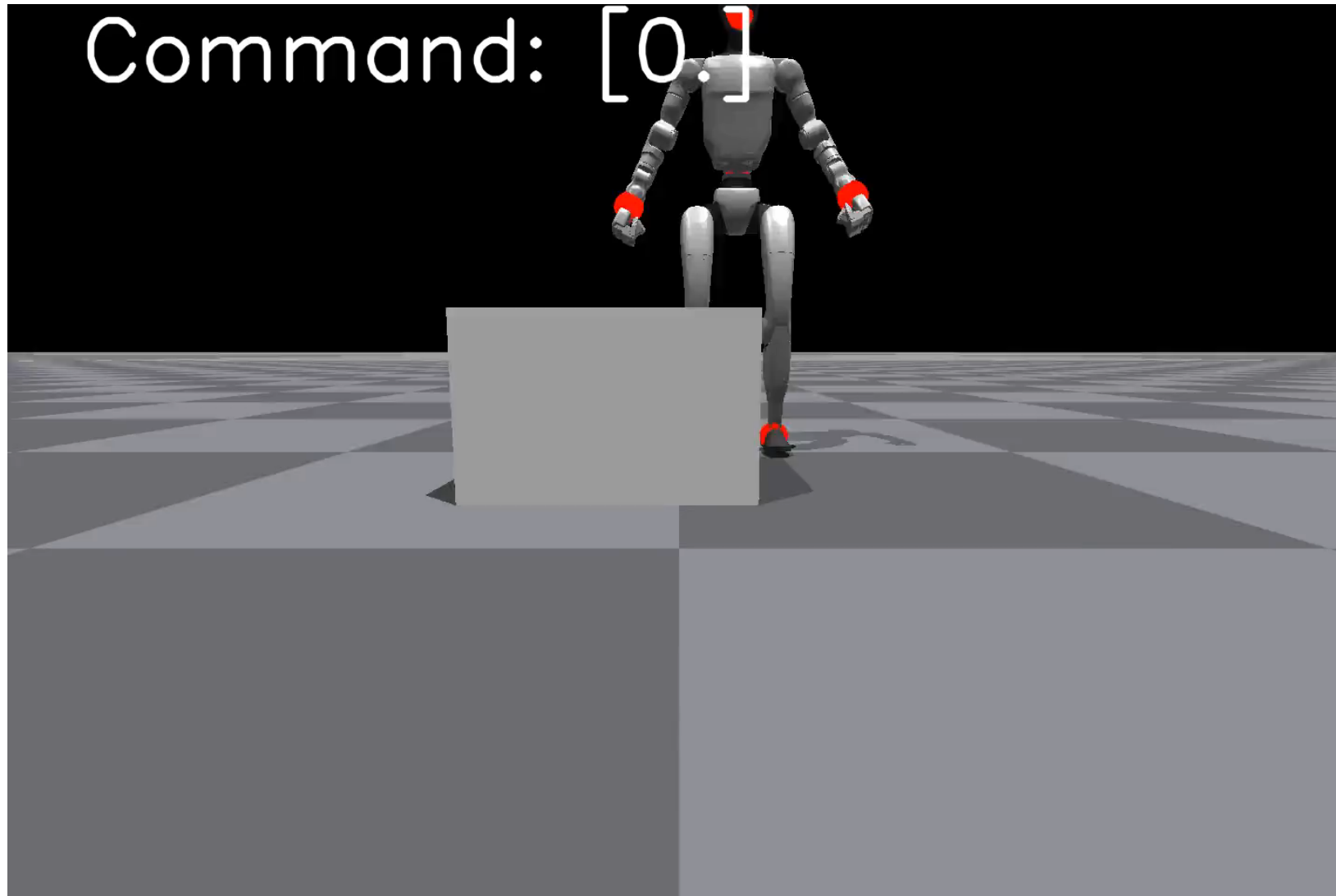


Real



- Six independently sampled rectangular masks
- Each mask covers **up to 25%** of the image area
- Each mask has a **10% probability** of being applied

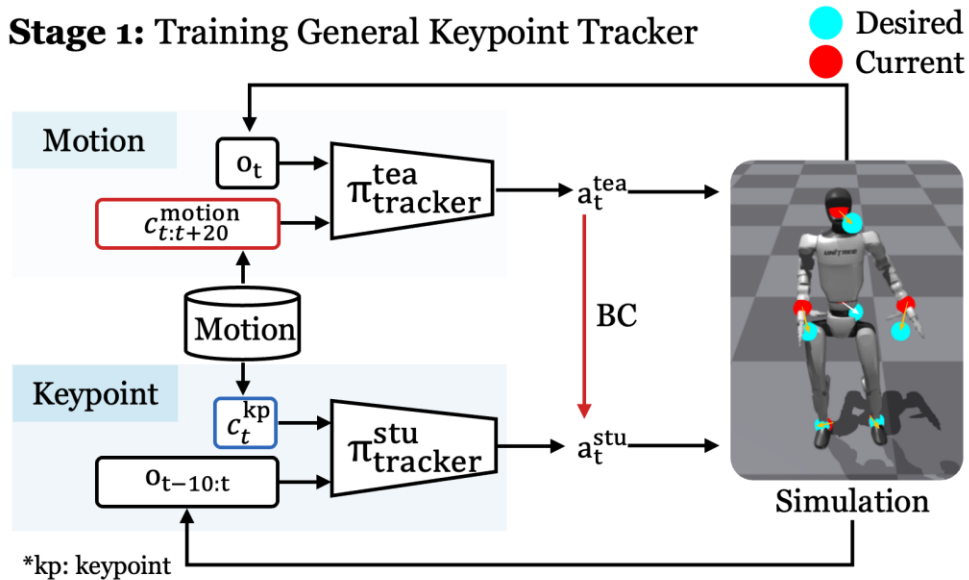
# Real World: Binary Command



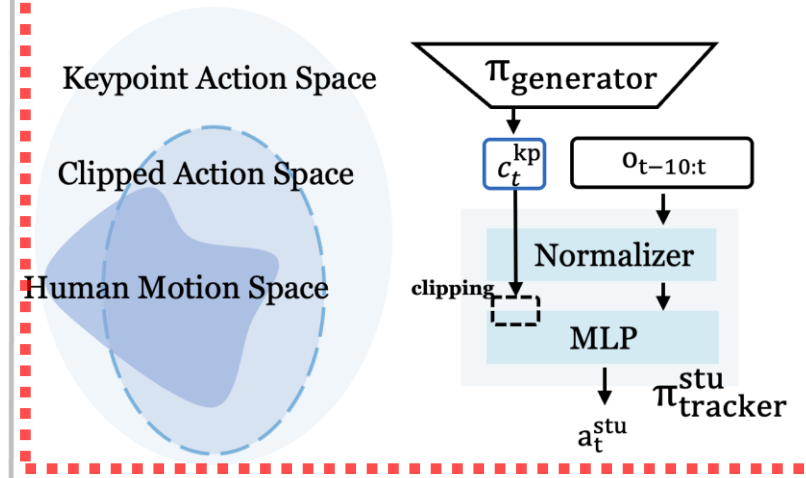


# Important Tricks 1 (Action Clip)

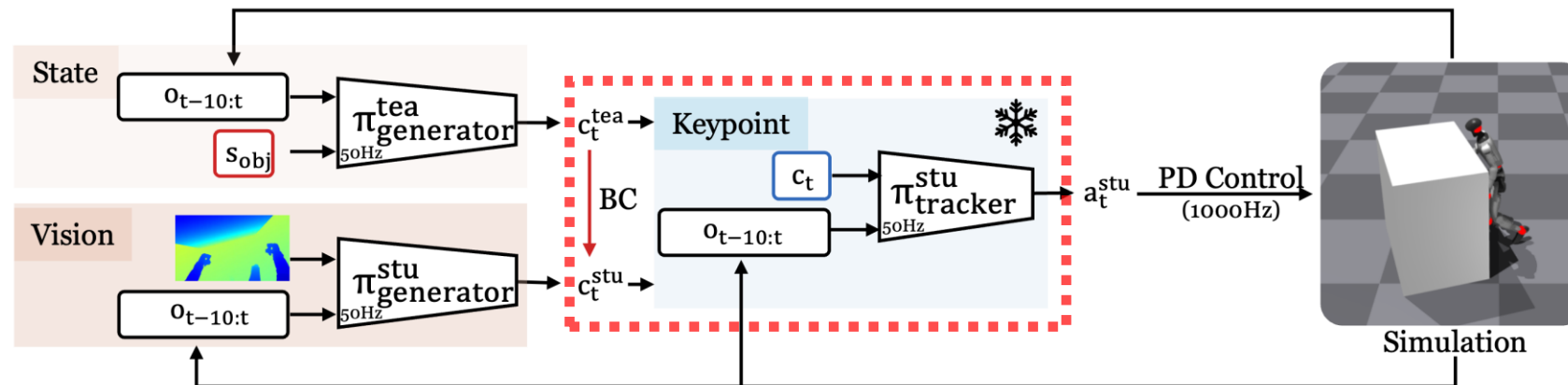
## Stage 1: Training General Keypoint Tracker



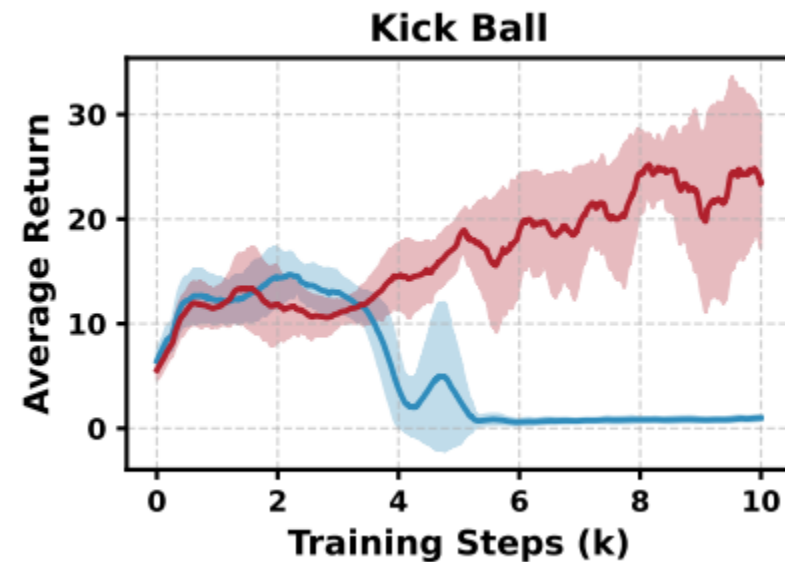
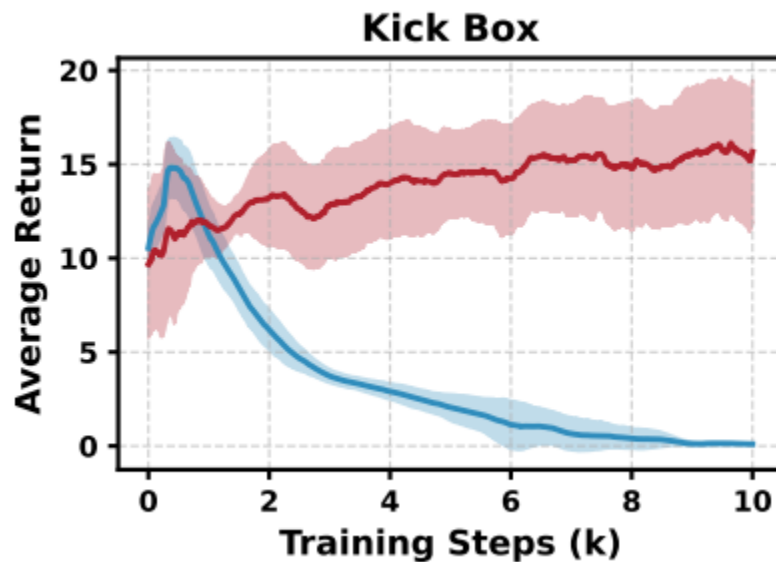
## Action Clipping via Human Motion Statistics



## Stage 2: Training Task-Specific Keypoint Generator



# Important Tricks 1 (Action Clip)



Method	Push Box			Kick Box			Kick Ball		
	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓
<b>State-based</b>									
Ours	<b>152 ± 36</b>	<b>151 ± 29</b>	13 ± 4	<b>78 ± 3</b>	<b>78 ± 3</b>	<b>0 ± 0</b>	<b>189 ± 3</b>	<b>189 ± 3</b>	4 ± 1
w/o clip	68 ± 118	67 ± 94	11 ± 16	3 ± 5	3 ± 4	<b>0 ± 0</b>	12 ± 15	12 ± 12	1 ± 1
<b>Vision-based</b>									
Ours	<b>37 ± 28</b>	<b>19 ± 15</b>	21 ± 12	<b>55 ± 5</b>	<b>30 ± 3</b>	33 ± 3	<b>135 ± 6</b>	<b>121 ± 8</b>	47 ± 12
w/o clip	10 ± 18	9 ± 12	4 ± 5	6 ± 7	5 ± 3	3 ± 3	1 ± 1	0 ± 1	<b>0 ± 0</b>

# Important Tricks 2 (Noise Augmentation)

## Implementation Details

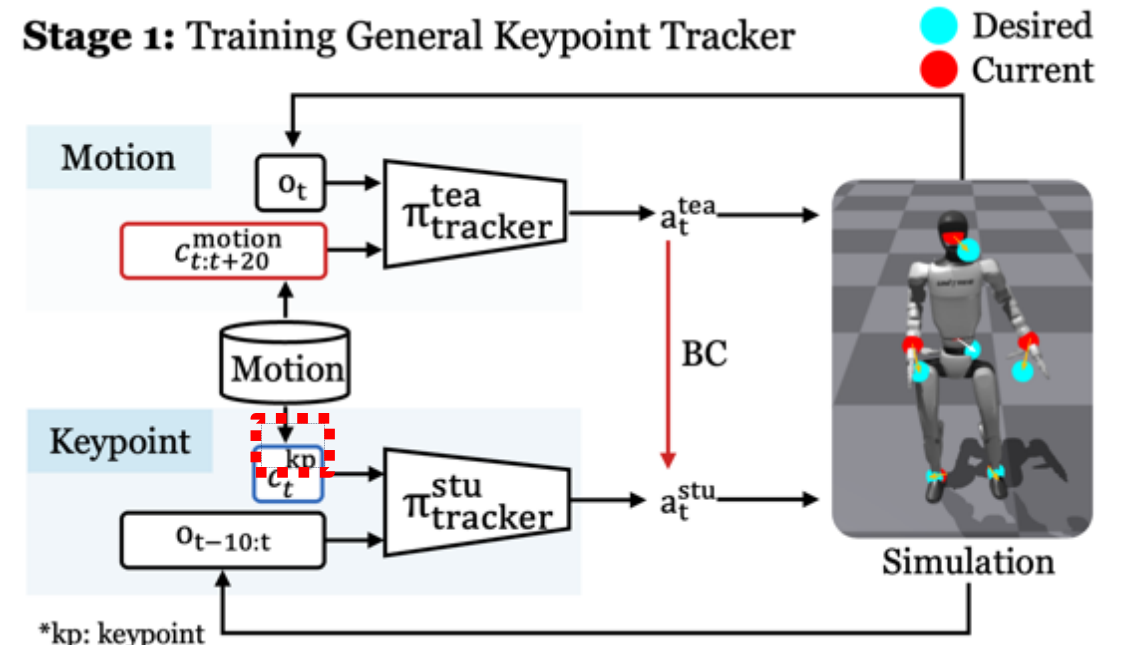
- Add noise for student input (Stage 1)

- relative noise

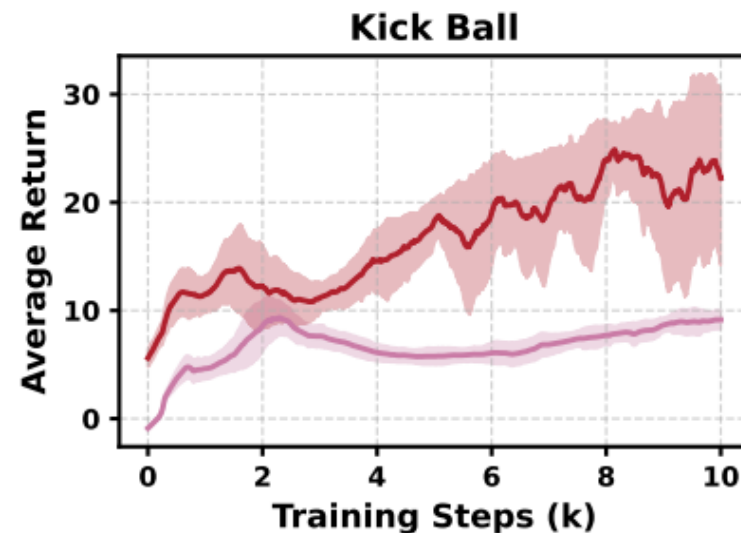
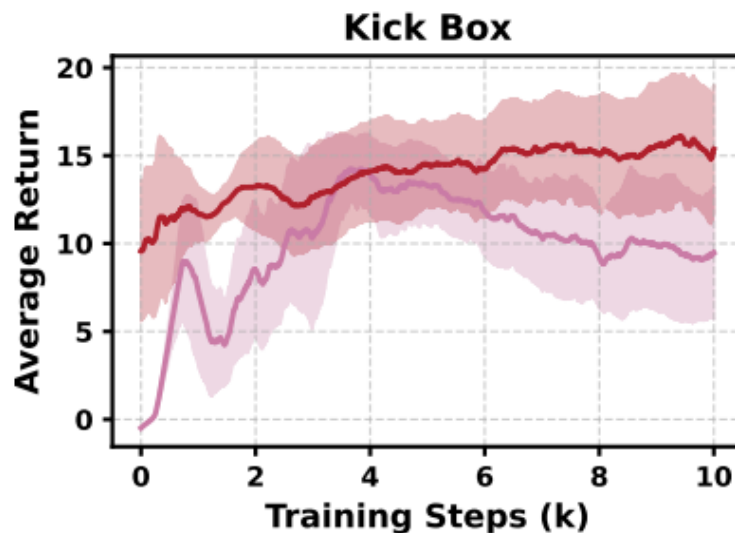
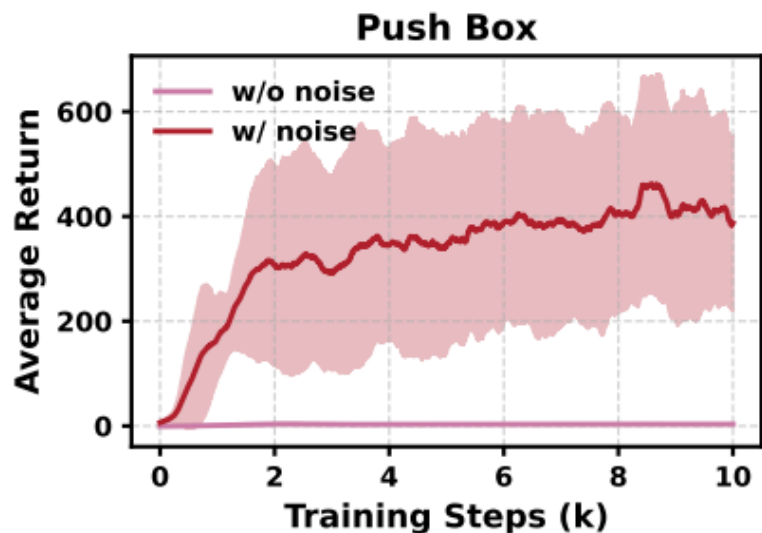
- $X_{\text{noised}} = X \cdot \lambda_i, \quad \forall i \in \{1, \dots, n\}, \quad \lambda_i \in [0.5, 1.5]$

- Why adding noise

- Kind of data augmentation



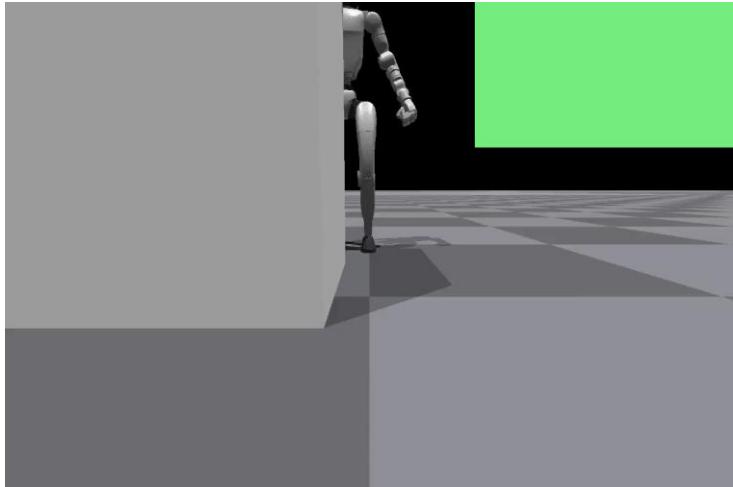
# Important Tricks 2 (Noise Augmentation)



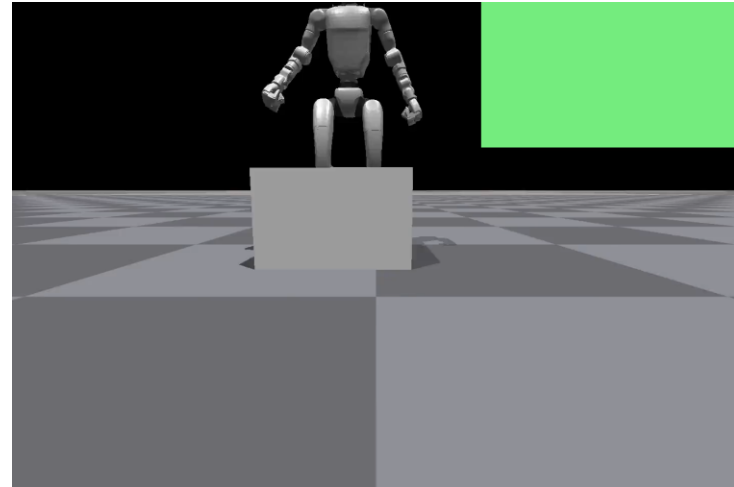
Method	Push Box			Kick Box			Kick Ball		
	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓
State-based									
Ours	<b>152 ± 36</b>	<b>151 ± 29</b>	13 ± 4	<b>78 ± 3</b>	<b>78 ± 3</b>	<b>0 ± 0</b>	<b>189 ± 3</b>	<b>189 ± 3</b>	4 ± 1
w/o noise	2 ± 1	2 ± 1	<b>0 ± 0</b>	30 ± 24	30 ± 20	1 ± 0	136 ± 8	136 ± 7	4 ± 0
Vision-based									
Ours	<b>37 ± 28</b>	<b>19 ± 15</b>	21 ± 12	<b>55 ± 5</b>	<b>30 ± 3</b>	33 ± 3	<b>135 ± 6</b>	<b>121 ± 8</b>	47 ± 12
w/o noise	2 ± 1	2 ± 1	<b>0 ± 0</b>	25 ± 7	11 ± 4	15 ± 3	86 ± 7	77 ± 7	30 ± 8



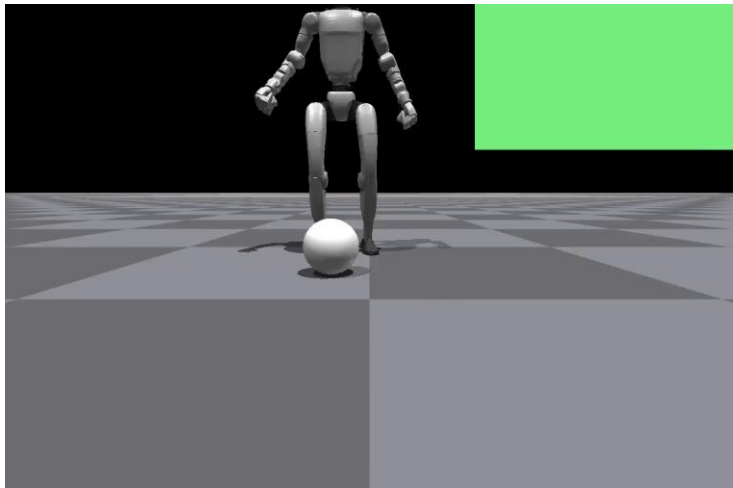
# Results



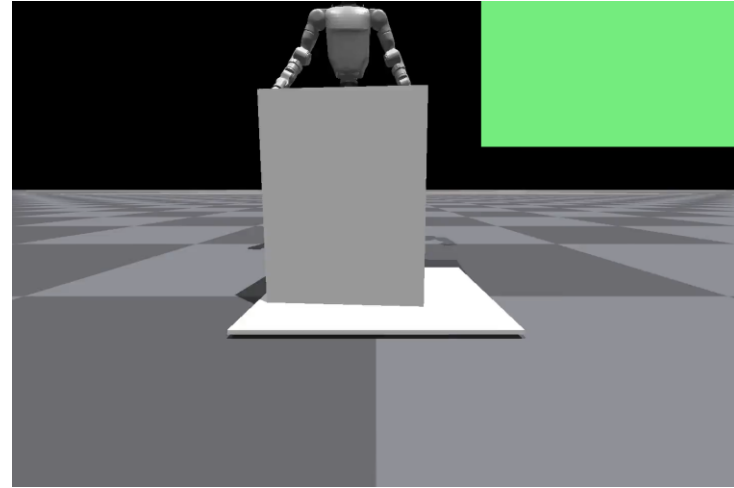
Push Box



Kick Box

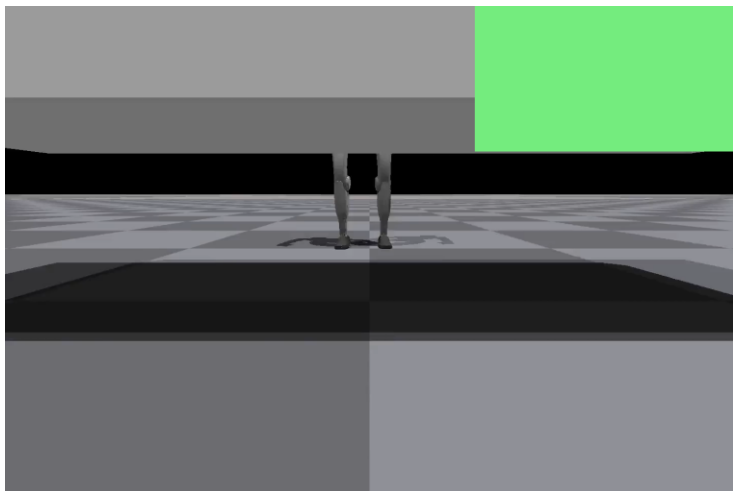


Kick Ball

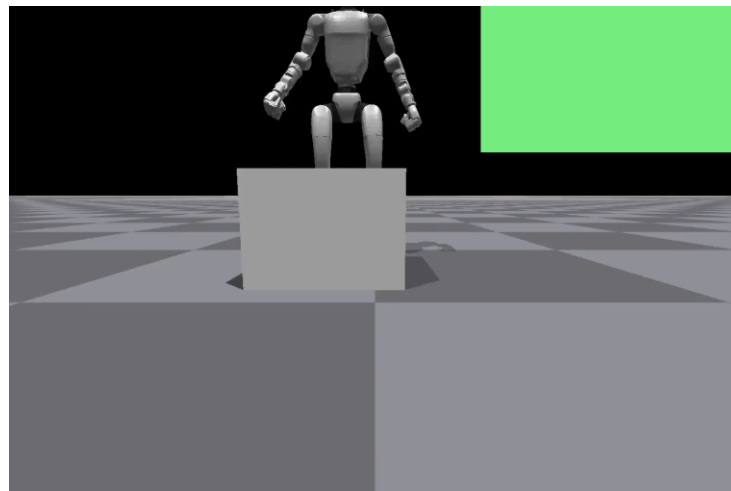


Lift Box

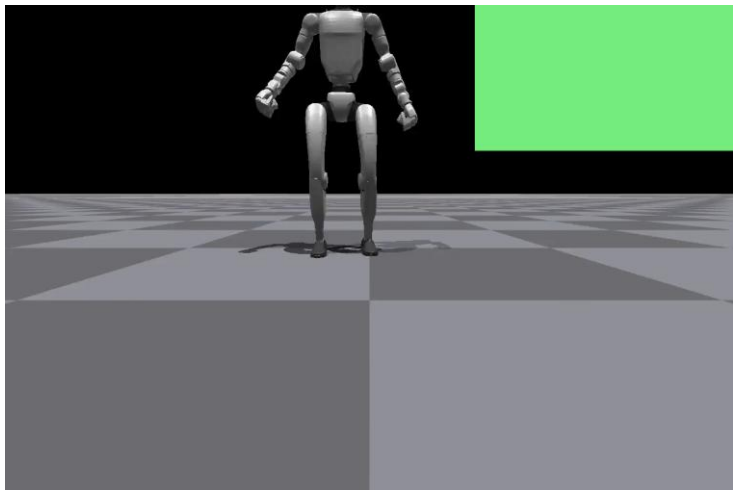
# Results



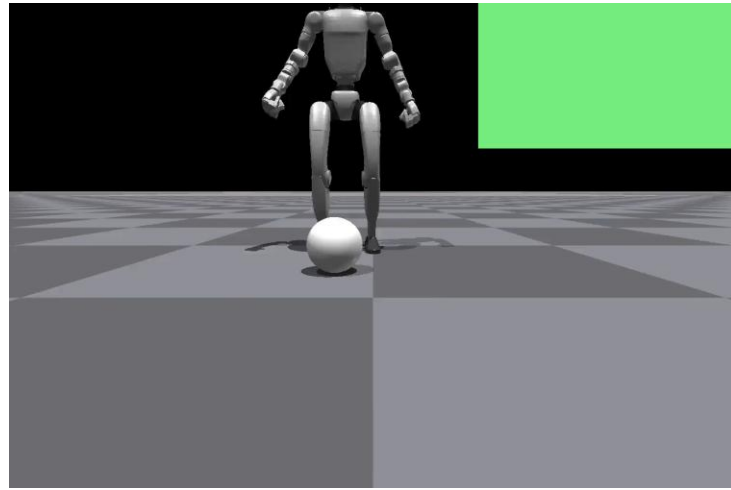
Push Cube



Large Kick



Reach Box

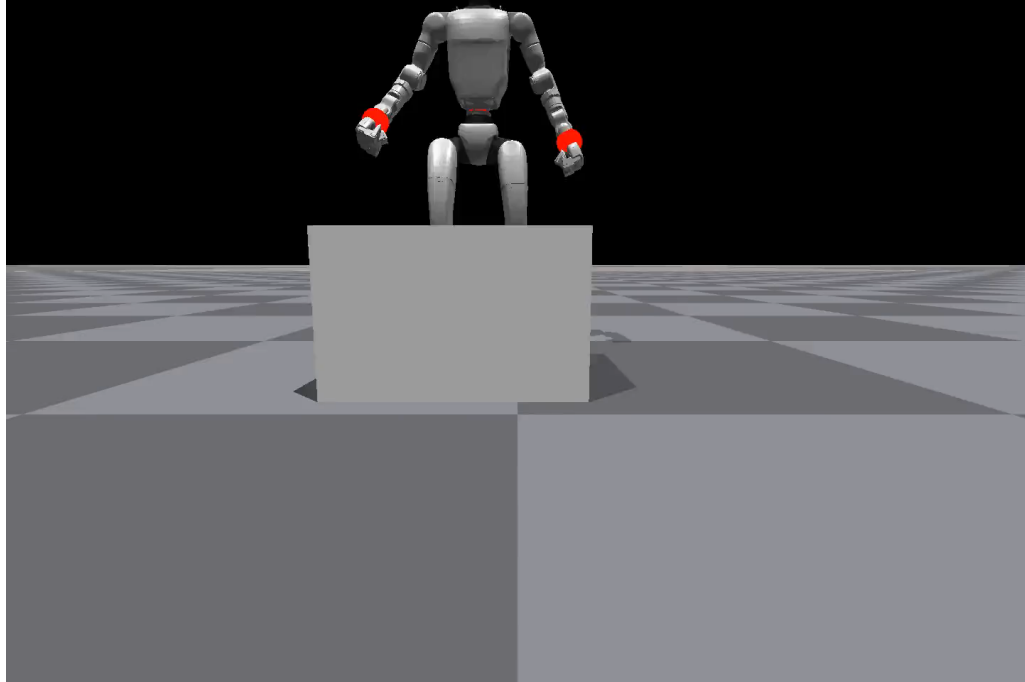


Balance Ball

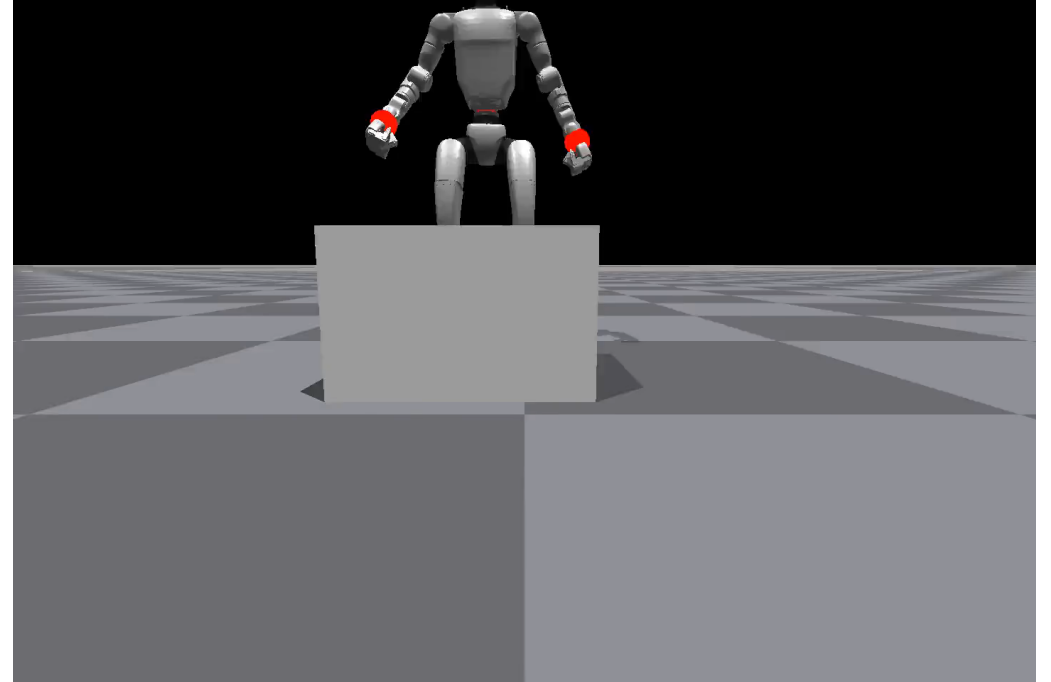
# Results (Simulation Evaluation)

Method	Push Box			Kick Box		
	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓
teacher	152 ± 36	151 ± 29	13 ± 4	78 ± 3	78 ± 3	0 ± 0
stu w/ vision	<b>37 ± 28</b>	<b>19 ± 15</b>	21 ± 12	<b>55 ± 5</b>	<b>30 ± 3</b>	33 ± 3
stu w/o vision	2 ± 0	2 ± 0	<b>1 ± 0</b>	0 ± 0	0 ± 0	<b>0 ± 0</b>
Method	Large Kick			Kick Ball		
	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓
teacher	8 ± 1	7 ± 1	2 ± 0	189 ± 3	189 ± 3	4 ± 1
stu w/ vision	<b>6 ± 0</b>	<b>6 ± 0</b>	2 ± 0	<b>135 ± 6</b>	<b>121 ± 8</b>	47 ± 12
stu w/o vision	4 ± 0	4 ± 0	<b>1 ± 0</b>	1 ± 0	1 ± 0	<b>0 ± 0</b>
Method	Lift Box			Reach Box		
	Height [m] ↑	Box Fall Rate [%] ↓	Alive [s] ↑	Velocity [m/s] ↑	Collision Rate [%] ↓	Alive [s] ↑
teacher	1 ± 0	34 ± 25	38 ± 13	4 ± 0	0 ± 0	60 ± 0
stu w/ vision	<b>1 ± 0</b>	30 ± 23	<b>30 ± 7</b>	<b>4 ± 0</b>	<b>0 ± 0</b>	<b>42 ± 6</b>
stu w/o vision	0 ± 0	<b>15 ± 21</b>	6 ± 4	<b>4 ± 0</b>	<b>0 ± 0</b>	18 ± 6
Method	Balance Ball			Push Cube (Tabletop)		
	Force [N] ↑	Foot Fall Rate [%] ↓	Alive [s] ↑	Error [cm] ↓	Finish Time [s] ↓	Alive [s] ↑
teacher	21 ± 2	0 ± 0	34 ± 8	9 ± 3	4 ± 1	58 ± 1
stu w/ vision	<b>24 ± 1</b>	<b>0 ± 0</b>	<b>45 ± 7</b>	<b>21 ± 2</b>	<b>20 ± 8</b>	<b>57 ± 0</b>
stu w/o vision	6 ± 0	<b>0 ± 0</b>	5 ± 1	57 ± 22	43 ± 8	51 ± 10

# Analysis Results (Training Pipeline)



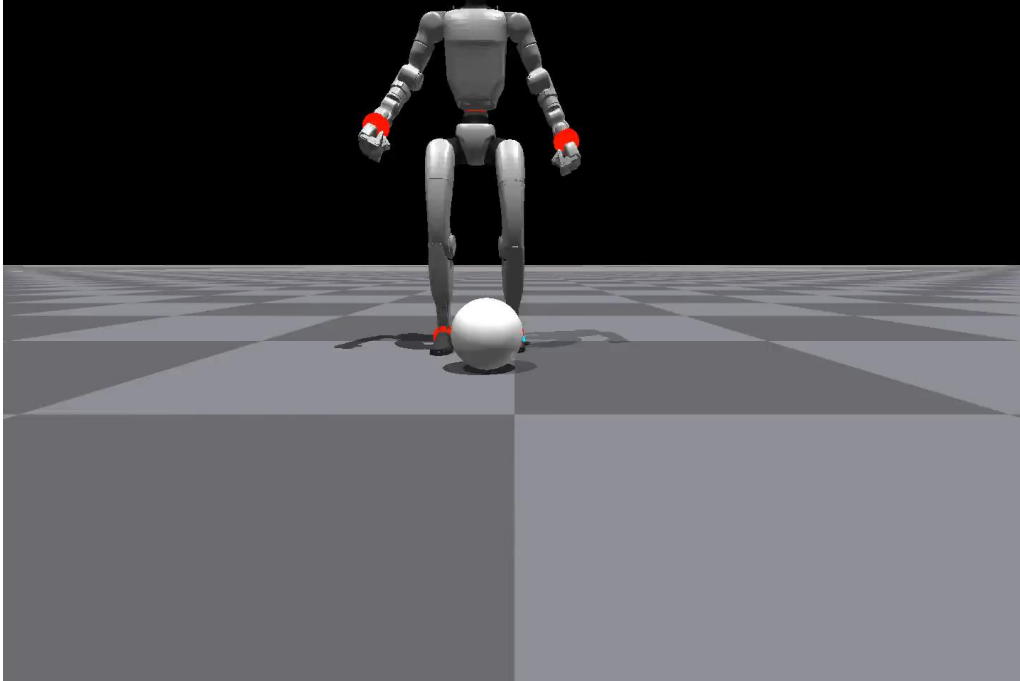
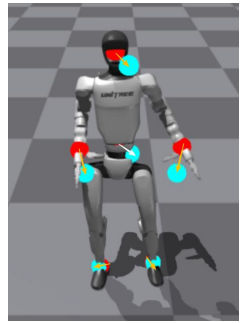
VisualMimic



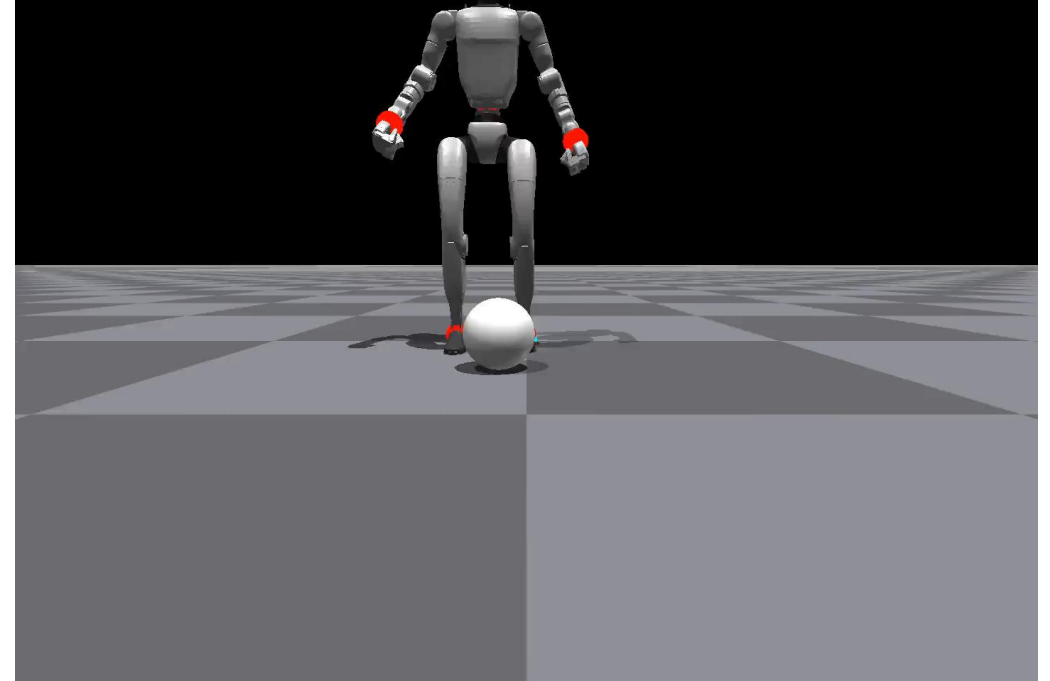
VisualMimic (w/o Stage 0&1 Distillation)



# Analysis Results (Interface Design)

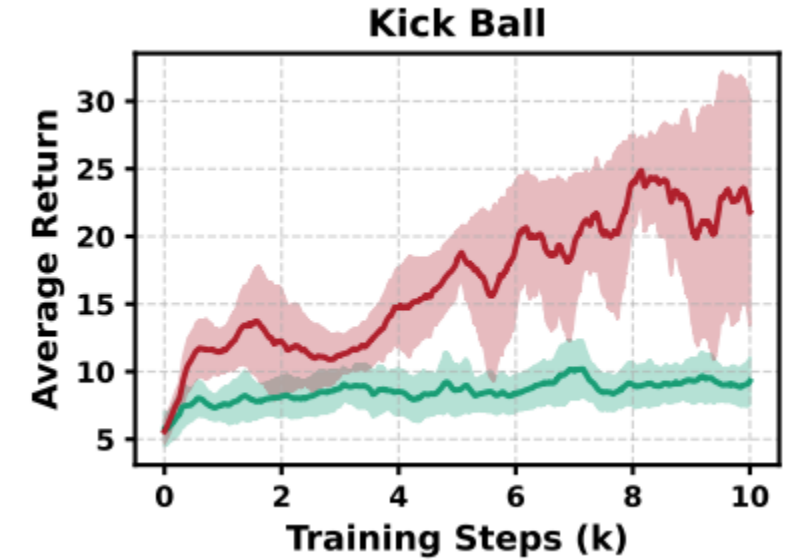
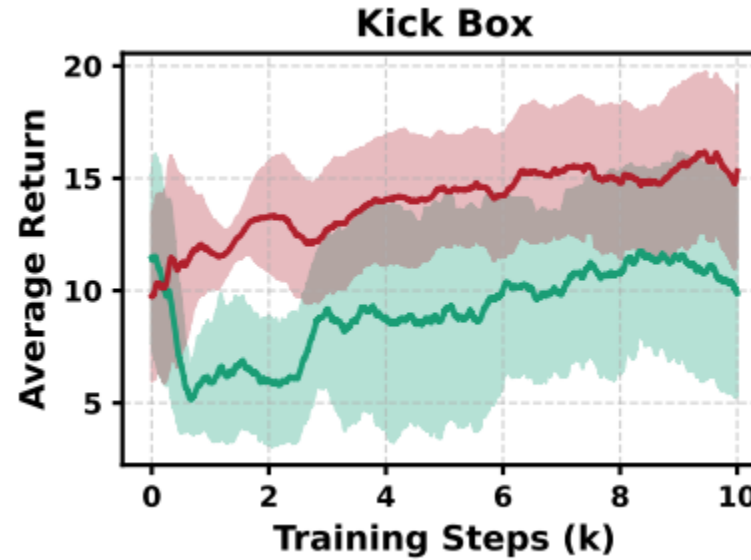
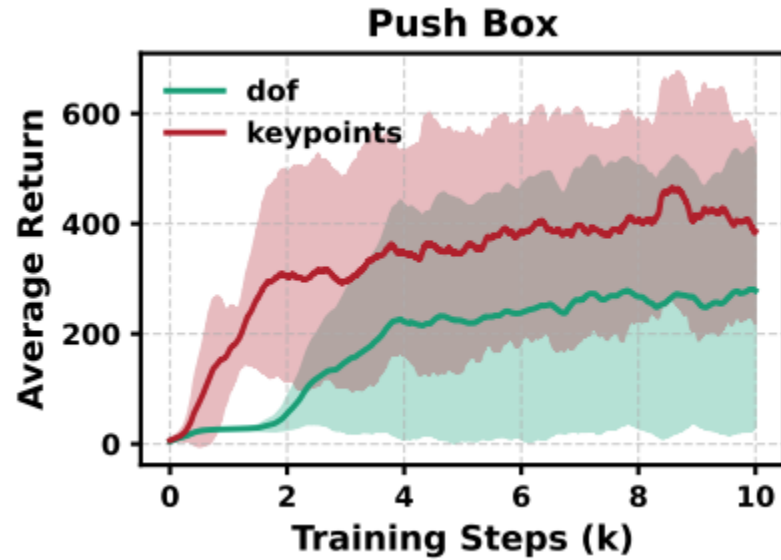


VisualMimic



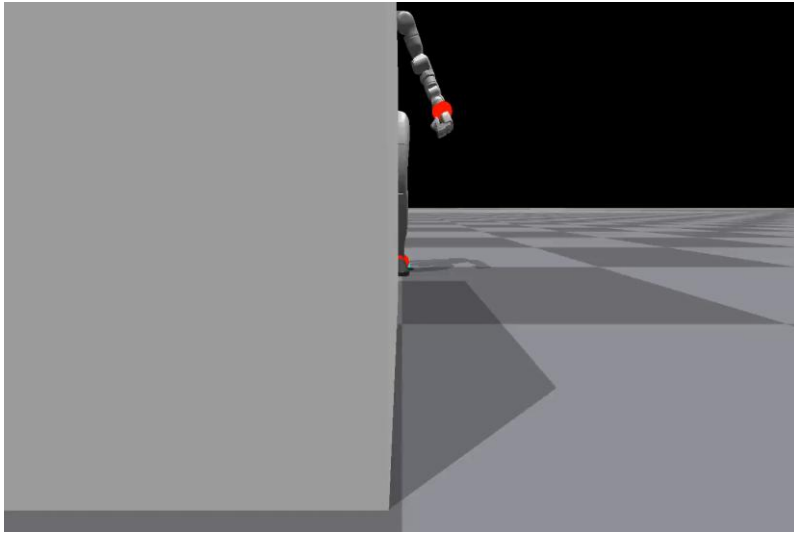
VisualMimic (w/ 3 point interface)

# Analysis Results (Interface Design)

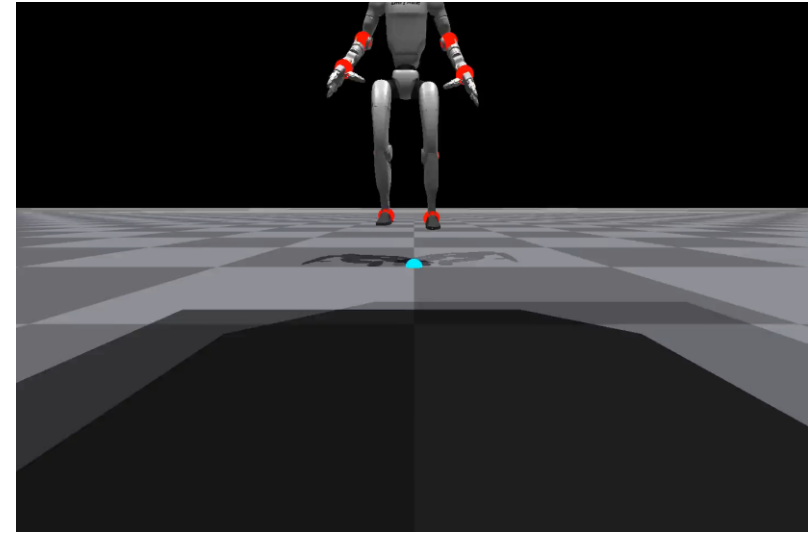


Method	Push Box			Kick Box			Kick Ball		
	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓
<b>State-based</b>									
Ours	<b>152 ± 36</b>	<b>151 ± 29</b>	13 ± 4	<b>78 ± 3</b>	<b>78 ± 3</b>	<b>0 ± 0</b>	<b>189 ± 3</b>	<b>189 ± 3</b>	4 ± 1
DoF as Interface	10 ± 9	8 ± 6	5 ± 3	40 ± 34	40 ± 28	<b>0 ± 0</b>	0 ± 0	0 ± 0	<b>0 ± 0</b>
<b>Vision-based</b>									
Ours	<b>37 ± 28</b>	<b>19 ± 15</b>	21 ± 12	<b>55 ± 5</b>	<b>30 ± 3</b>	33 ± 3	<b>135 ± 6</b>	<b>121 ± 8</b>	47 ± 12
DoF as Interface	10 ± 2	6 ± 1	6 ± 1	5 ± 4	1 ± 0	4 ± 3	0 ± 0	0 ± 0	<b>0 ± 0</b>

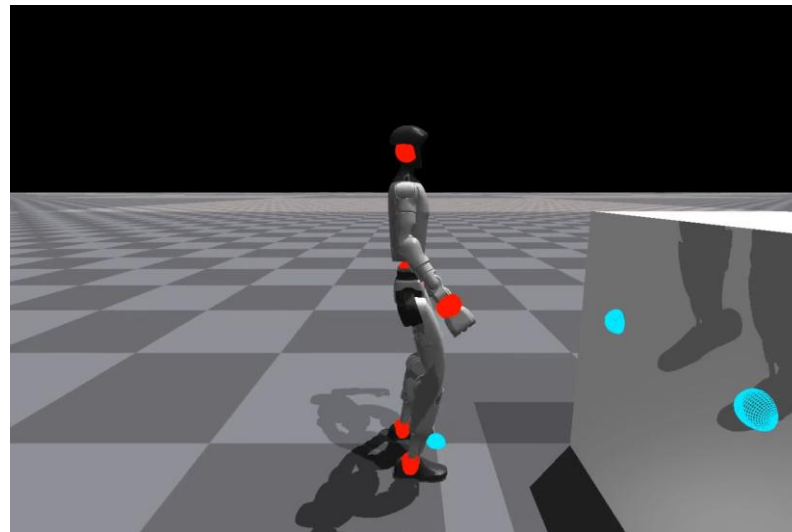
# Analysis Results (Whole-body Dexterity)



Foot and Hand



Shoulder



Two Hands

# Analysis Results

Method	Push Box			Kick Box			Kick Ball		
	Distance [m] $\uparrow$	Forward [m] $\uparrow$	Drift [m] $\downarrow$	Distance [m] $\uparrow$	Forward [m] $\uparrow$	Drift [m] $\downarrow$	Distance [m] $\uparrow$	Forward [m] $\uparrow$	Drift [m] $\downarrow$
<b>State-based</b>									
Ours	<b>152 <math>\pm</math> 36</b>	<b>151 <math>\pm</math> 29</b>	13 $\pm$ 4	<b>78 <math>\pm</math> 3</b>	<b>78 <math>\pm</math> 3</b>	<b>0 <math>\pm</math> 0</b>	<b>189 <math>\pm</math> 3</b>	<b>189 <math>\pm</math> 3</b>	4 $\pm$ 1
w/o noise	2 $\pm$ 1	2 $\pm$ 1	<b>0 <math>\pm</math> 0</b>	30 $\pm$ 24	30 $\pm$ 20	1 $\pm$ 0	136 $\pm$ 8	136 $\pm$ 7	4 $\pm$ 0
w/o clip	68 $\pm$ 118	67 $\pm$ 94	11 $\pm$ 16	3 $\pm$ 5	3 $\pm$ 4	<b>0 <math>\pm</math> 0</b>	12 $\pm$ 15	12 $\pm$ 12	1 $\pm$ 1
DoF as Interface	10 $\pm$ 9	8 $\pm$ 6	5 $\pm$ 3	40 $\pm$ 34	40 $\pm$ 28	<b>0 <math>\pm</math> 0</b>	0 $\pm$ 0	0 $\pm$ 0	<b>0 <math>\pm</math> 0</b>
Local-Frame Tracker	38 $\pm$ 27	30 $\pm$ 23	16 $\pm$ 15	45 $\pm$ 7	45 $\pm$ 5	1 $\pm$ 0	109 $\pm$ 23	109 $\pm$ 19	7 $\pm$ 1
<b>Vision-based</b>									
Ours	<b>37 <math>\pm</math> 28</b>	<b>19 <math>\pm</math> 15</b>	21 $\pm$ 12	<b>55 <math>\pm</math> 5</b>	<b>30 <math>\pm</math> 3</b>	33 $\pm$ 3	<b>135 <math>\pm</math> 6</b>	<b>121 <math>\pm</math> 8</b>	47 $\pm$ 12
w/o noise	2 $\pm$ 1	2 $\pm$ 1	<b>0 <math>\pm</math> 0</b>	25 $\pm$ 7	11 $\pm$ 4	15 $\pm$ 3	86 $\pm$ 7	77 $\pm$ 7	30 $\pm$ 8
w/o clip	10 $\pm$ 18	9 $\pm$ 12	4 $\pm$ 5	6 $\pm$ 7	5 $\pm$ 3	3 $\pm$ 3	1 $\pm$ 1	0 $\pm$ 1	<b>0 <math>\pm</math> 0</b>
DoF as Interface	10 $\pm$ 2	6 $\pm$ 1	6 $\pm$ 1	5 $\pm$ 4	1 $\pm$ 0	4 $\pm$ 3	0 $\pm$ 0	0 $\pm$ 0	<b>0 <math>\pm</math> 0</b>
Local-Frame Tracker	14 $\pm$ 11	7 $\pm$ 5	8 $\pm$ 6	27 $\pm$ 15	16 $\pm$ 9	15 $\pm$ 6	38 $\pm$ 17	34 $\pm$ 13	12 $\pm$ 4
Visual RL	25 $\pm$ 16	11 $\pm$ 6	16 $\pm$ 9	0 $\pm$ 0	0 $\pm$ 0	<b>0 <math>\pm</math> 0</b>	0 $\pm$ 0	0 $\pm$ 0	<b>0 <math>\pm</math> 0</b>
<b>Blind</b>									
Ours w/o vision	2 $\pm$ 0	2 $\pm$ 0	1 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	1 $\pm$ 0	1 $\pm$ 0	0 $\pm$ 0



# Takeaway & Limitation

- RL is a beast
  - Powerful
  - Without proper DR or termination, it might bite
- Human motion matters
- Surprising facts
  - The visual sim2real gap isn't as large as expected.
  - Distilled visual policies still lag far behind their state-based counterparts.

<https://visualmimic.github.io>

Code, Video

# Thank You!

Contact: [yinshaofeng04@gmail.com](mailto:yinshaofeng04@gmail.com)

Website: <https://operator22th.github.io/>

I am applying for PhD programs in Fall 2026.

I'd love to connect and discuss research opportunities.