

# Deep Whole-Body Control

## Learning a Unified Policy for Manipulation and Locomotion

Zipeng Fu\* Xuxin Cheng\* Deepak Pathak

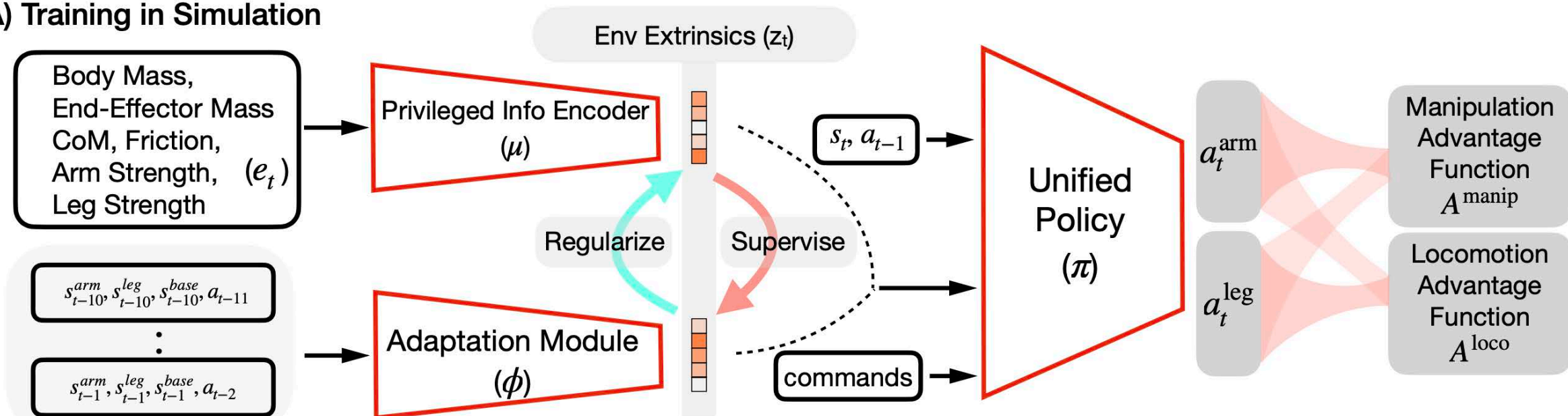


Videos  
&  
Code

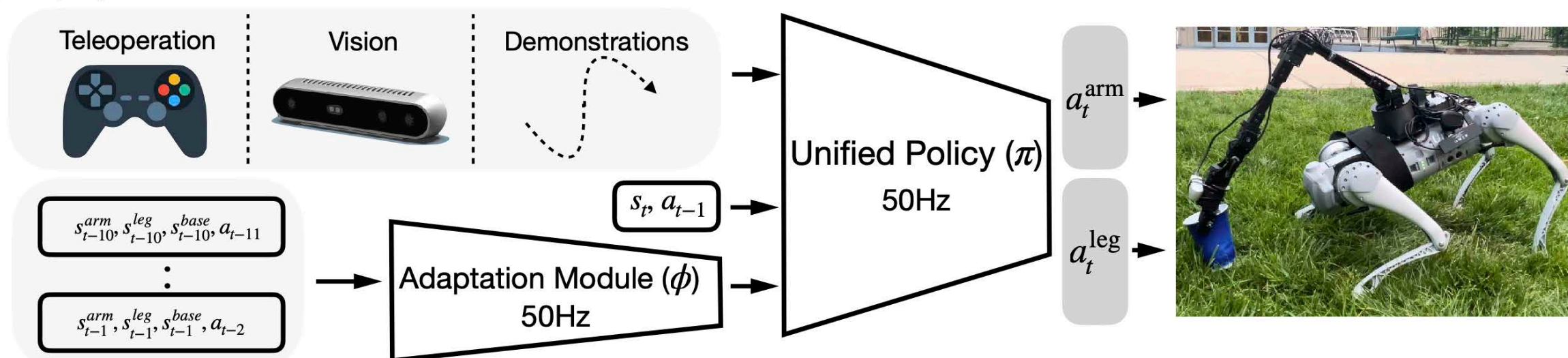
<https://maniploco.github.io>

## DeepWBC Pipeline

### (A) Training in Simulation



### (B) Deployment in Real World



**Motivation:** whole-body >> modular, low-cost >> expensive hardware

**TL;DR:** learning an **end-to-end** unified policy for whole-body control of a custom-built **low-cost** quadruped **mobile manipulator**

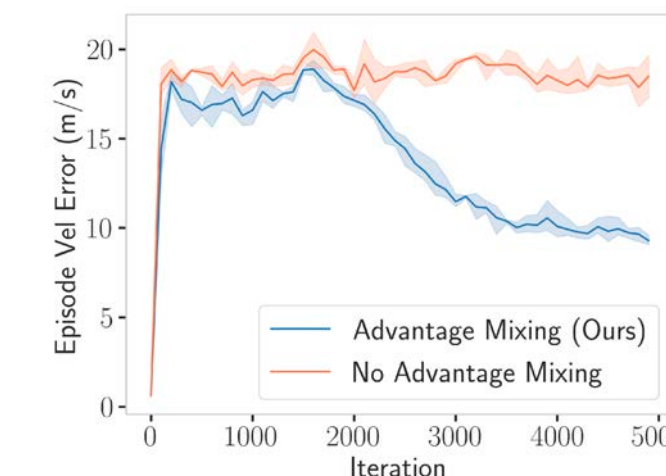
**Regularized Online Adaptation for Sim-to-Real Transfer** ( $\lambda, \beta$  follow simple linear curricula)

$$L(\theta_\pi, \theta_\mu, \theta_\phi) = \underbrace{-J(\theta_\pi, \theta_\mu)}_{\text{RL Loss}} + \underbrace{\lambda \|z^\mu - \text{sg}[z^\phi]\|_2}_{\text{Regularization}} + \underbrace{\|\text{sg}[z^\mu] - z^\phi\|_2}_{\text{Adaptation}}$$

**Advantage Mixing for Policy Learning**

$$J(\theta_\pi) = \frac{1}{|\mathcal{D}|} \sum_{(s_t, a_t) \in \mathcal{D}} \log \pi(a_t^{\text{arm}} | s_t) (A^{\text{manip}} + \beta A^{\text{loco}}) + \log \pi(a_t^{\text{leg}} | s_t) (\beta A^{\text{manip}} + A^{\text{loco}})$$

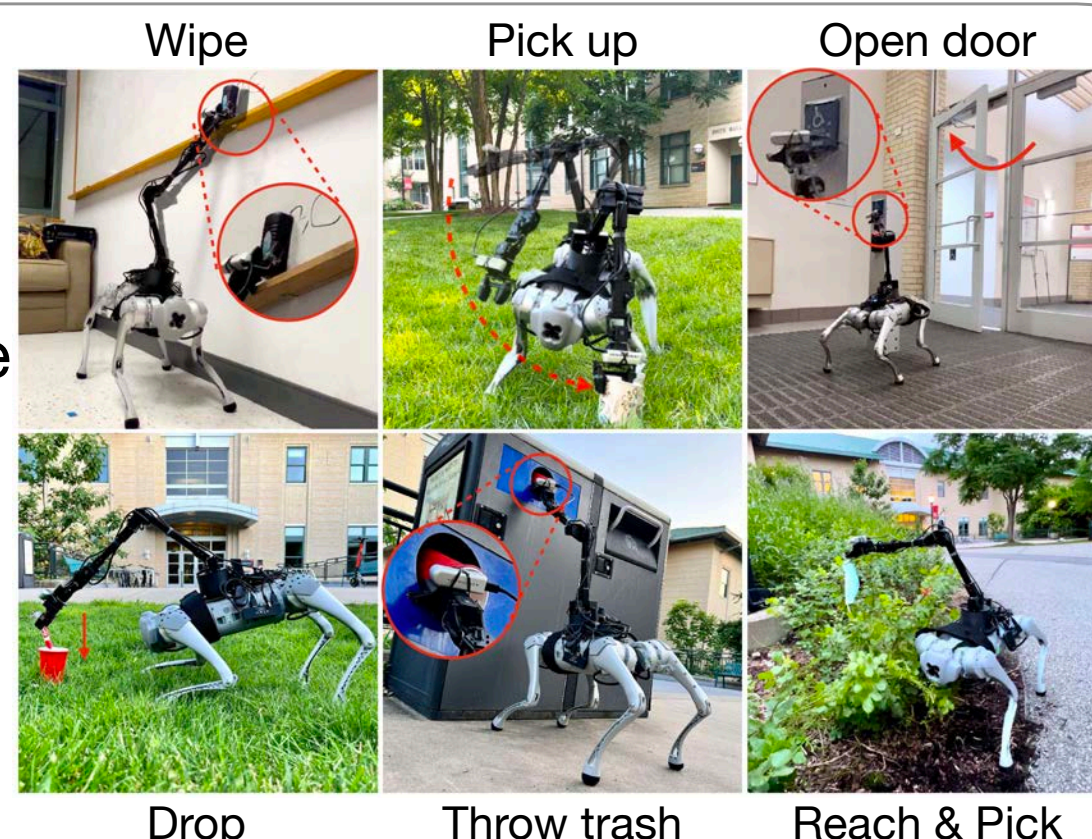
	Realizability Gap $\ z^\mu - z^\phi\ _2 \downarrow$	Survival $\uparrow$	Vel Error $\downarrow$	EE Error $\downarrow$
Domain Randomization	-	$95.8 \pm 0.2$	$0.46 \pm 0.00$	$0.40 \pm 0.00$
RMA [22]	$0.31 \pm 0.01$	$95.2 \pm 0.2$	$0.44 \pm 0.00$	$0.26 \pm 0.04$
Regularized Online Adapt (Ours)	<b><math>2\text{e-}4 \pm 0.00</math></b>	<b><math>97.4 \pm 0.1</math></b>	<b><math>0.39 \pm 0.01</math></b>	<b><math>0.21 \pm 0.00</math></b>
Expert w/ Reg.	-	$97.8 \pm 0.2$	$0.40 \pm 0.01$	$0.21 \pm 0.00$
Expert w/o Reg.	-	$98.3 \pm 0.2$	$0.39 \pm 0.00$	$0.21 \pm 0.00$



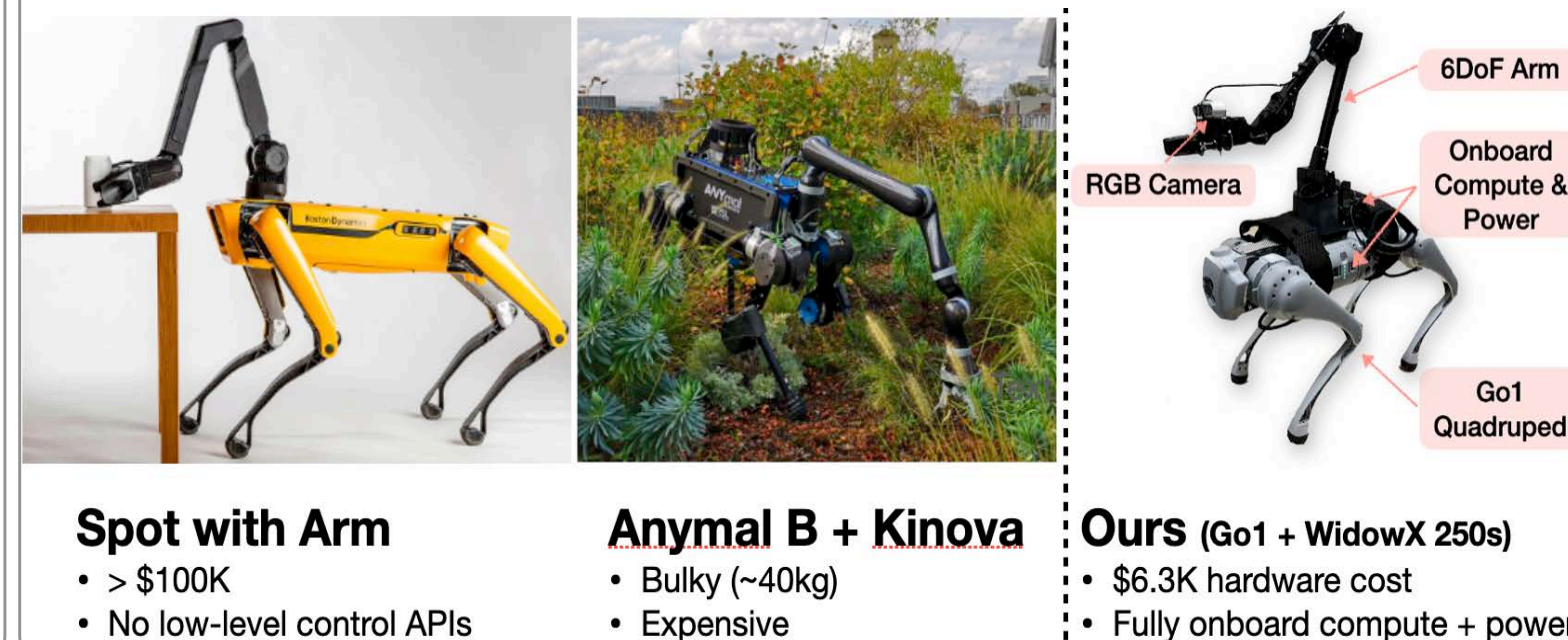
## Large Whole-Body Workspace



## Diverse Task Set



## Hardware Setup Comparison



## Ours

## Modular Baseline (MPC+IK)

