

# From **Human**, Go Beyond: The Robotic Era

Shanghai Qizhi Institute & Shanghai Jiao Tong University

Kun Lei

13.01.2026



# What the robot needs to do ?



- Reliability
- Efficiency
- Robustness

## □ Motivation

- The pre-training and post-training of our human beings



“**Pre-training**” under supervision of parents



“Self-supervised **finetuning**”

Improve **generalization** ability ?

- We finetune ourself for a lot of objectives, such as ...



acceleration



precision



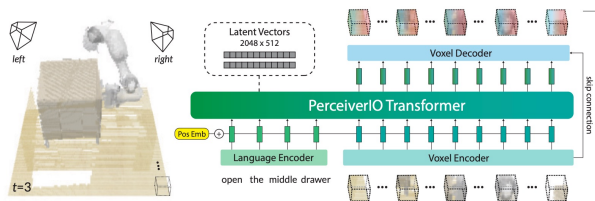
Robustness & safety



# □ Motivation

- Learn from demonstration, but also limited by the dataset

## PerACT



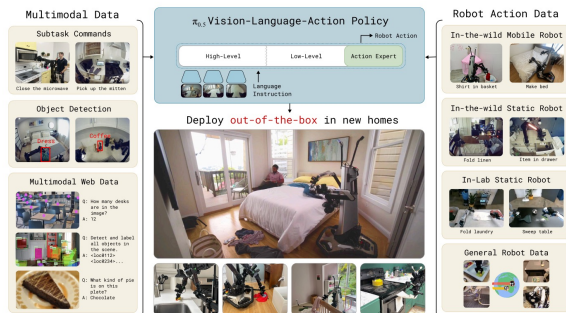
$\pi_0$ : A Vision-Language-Action Flow Model for General Robot Control

**Physical Intelligence**  
Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolò Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith Mohanakrishnan, Suraj Nair, Karl Pertsch, Lucy Xuyang Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, Ury Zhilinsky  
<https://physicalintelligence.company/blog/pi0>



$\pi_{0.5}$ : a Vision-Language-Action Model with Open-World Generalization

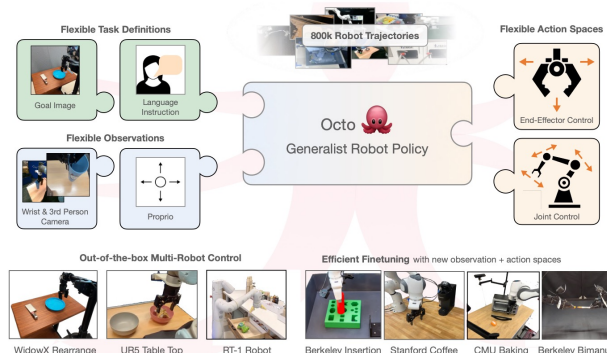
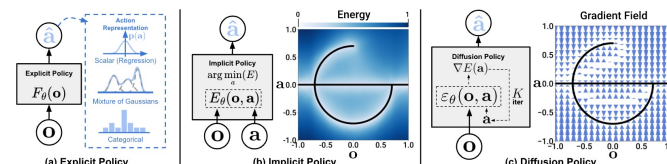
**Physical Intelligence**  
Kevin Black, Noah Brown, James Darpinian, Karam Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolò Fusai, Manuel Y. Galkner, Dibya Ghosh, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mohanakrishnan, Suraj Nair, Karl Pertsch, Allen Z. Ren, Lucy Xuyang Shi, Laura Smith, Just Tobias Springenberg, Kyle Stachowicz, James Tanner, Quan Vuong, Homer Walker, Anna Walling, Haohuan Wang, Lili Yu, Ury Zhilinsky  
<https://physicalintelligence.company/blog/pi0.5>



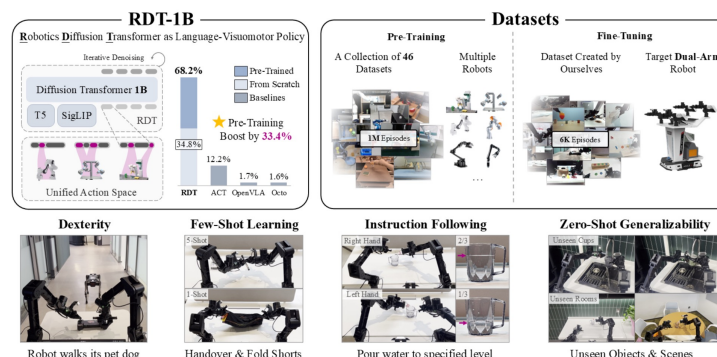
$\pi_{0.5}$

## Diffusion Policy

Visuomotor Policy Learning via Action Diffusion



## Octo

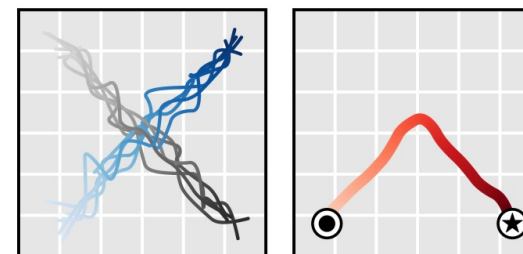


RDT

- ✓ Imitation Learning
- Simple and efficient
- 😭 Drawbacks:
  - Expert demonstration
  - Poor generalization ability



- ✓ Offline RL
- Learn from Suboptimal data
- Better generalization ability



“Stich”

What can we do if the era of “GPT-2” in robotics really comes ?

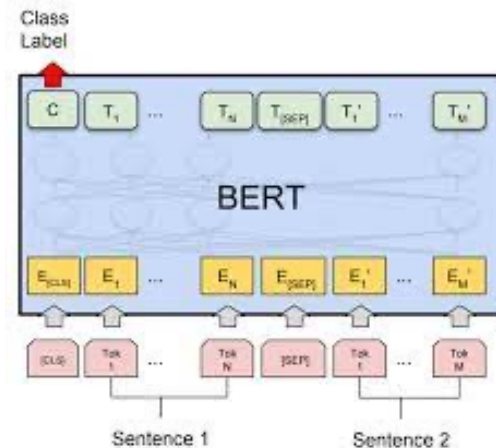


## □ Motivation

- Supervised Learning is great, but also limited by the dataset

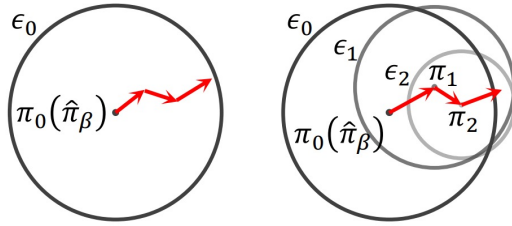


- How this paradigm work for robotic?
  - ✓ Imitation learning
  - ✓ Offline reinforcement learning
  - ✓ Offline to online RL finetuning



Successes of this paradigm in research areas of CV and NLP

## □ What we did toward this goal



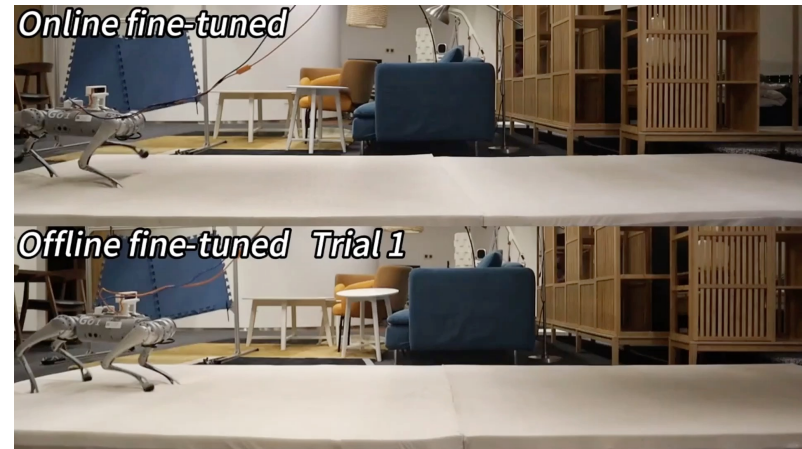
### 1. Behavior proximal policy optimization

Zifeng Zhuang\*, **Kun Lei\***, Jinxin Liu, Donglin Wang, Yilang Guo.  
*International Conference on Learning Representations (ICLR), 2023*

### 2. Uni-O4: Unifying Online and Offline Deep Reinforcement Learning with Multi-Step On-Policy Optimization

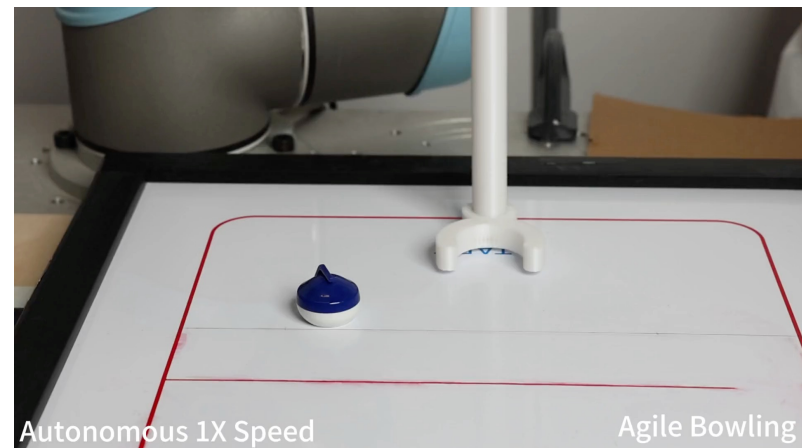
**Kun Lei**, Zhengmao He\*, Chenhao Lu\*, Kaizhe Hu,  
Yang Gao, Huazhe Xu.

*International Conference on Learning Representations (ICLR), 2024*



### 3. RL-100: Performant Robotic Manipulation with Real-World Reinforcement Learning

**Kun Lei\***, Huanyu Li\*, Dongjie Yu\*, Zhenyu Wei\*, Lingxiao Guo,  
Zhennan Jiang, Ziyu Wang, Shiyu Liang, Huazhe Xu.  
Preparing to submit.



# The post-training of robot (foundation) model

**Offline reinforcement learning**  
**- Pretraining/Finetuning**

Offline-to-online finetuning

- for fast adaptation
- with safety consideration
- for task acceleration

Finetuning from multi-modal perception



# Behavior proximal policy optimization

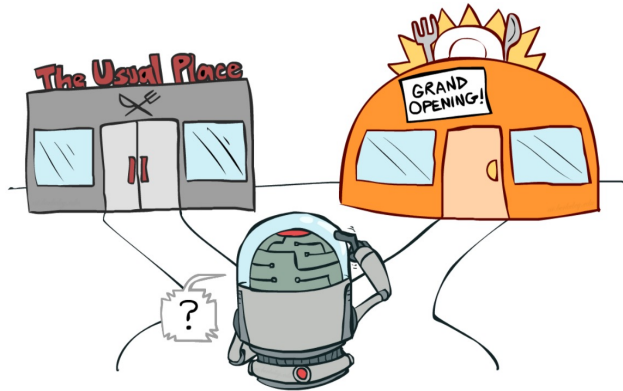
Zifeng Zhuang\*, **Kun Lei\***, Jinxin Liu, Donglin Wang, Yilang Guo



# ❑ Behavior proximal policy optimization (BPPO)

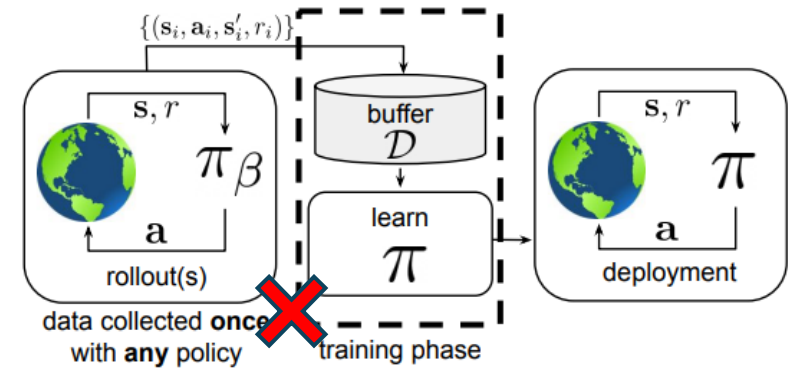
## ❑ Online RL

✓ Exploration is crucial

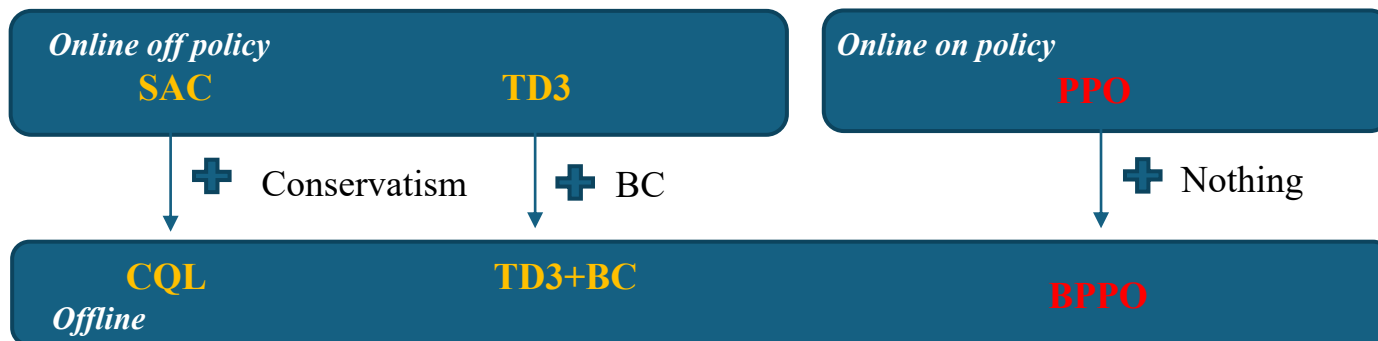


## ❑ Offline RL

😞 Exploration is limited



## ❑ Conservative methods



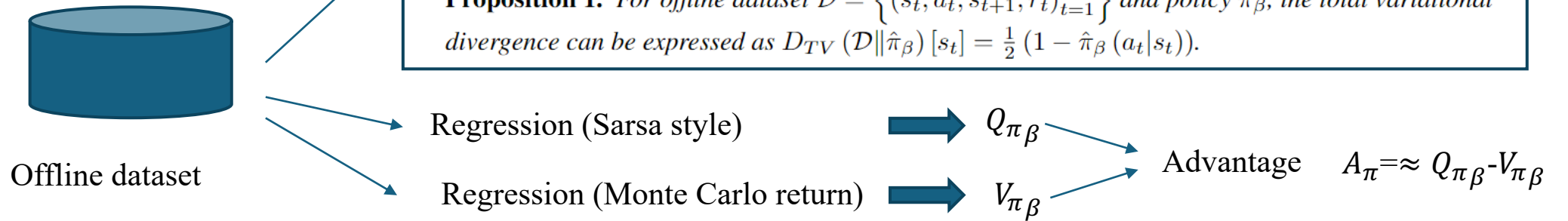
“An offline version of PPO”

Key insight:

The first work paid attention to **the policy learning** instead of the **value learning** in offline setting

## □ Method

### ➤ Stage 1: Supervised Learning



### ➤ Stage 2: Policy improvement using PPO

## □ Offline monotonic improvement over behavior policy

- For two policies  $\pi$  and  $\pi'$ , the **Performance Difference**  $J_\Delta(\pi', \pi)$  can be measured by the advantage function:

Online:  $J_\Delta(\pi', \pi) = \mathbb{E}_{\tau \sim P_{\pi'}(\tau)} \left[ \sum_{t=0}^H \gamma^t A_\pi(s_t, a_t) \right] = \mathbb{E}_{s \sim \rho_{\pi'}(\cdot), a \sim \pi'(\cdot | s)} [A_\pi(s, a)]$

Offline:  $\hat{J}_\Delta(\pi, \hat{\pi}_\beta) = \mathbb{E}_{s \sim \rho_{\mathcal{D}}(\cdot), a \sim \pi(\cdot | s)} [A_{\hat{\pi}_\beta}(s, a)]$

From TRPO

**Theorem 2.** Given the distance  $D_{TV}(\pi \parallel \hat{\pi}_\beta)[s]$  and  $D_{TV}(\mathcal{D} \parallel \hat{\pi}_\beta)[s] = \frac{1}{2} (1 - \hat{\pi}_\beta(a_t | s_t))$ , we can

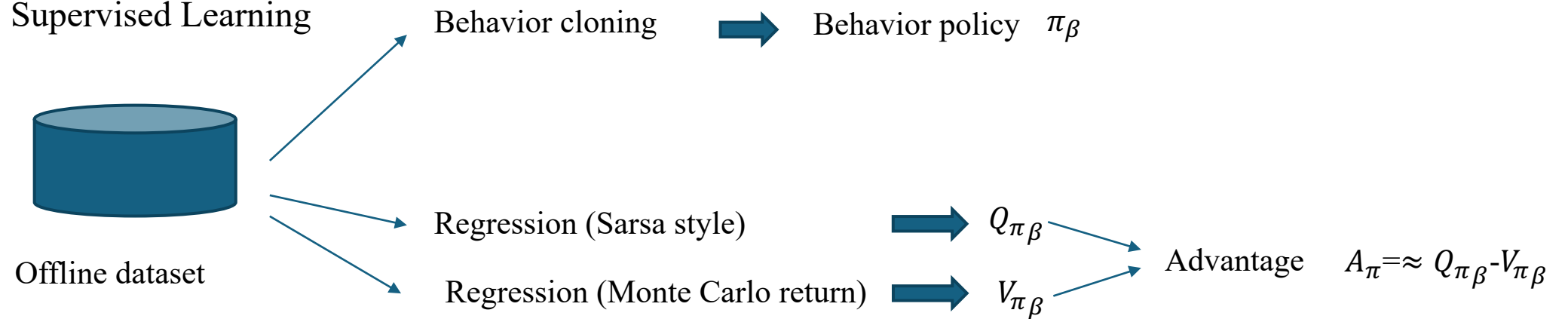
Only in offline

$$J_\Delta(\pi, \hat{\pi}_\beta) \geq \hat{J}_\Delta(\pi, \hat{\pi}_\beta) - 4\gamma \mathbb{A}_{\hat{\pi}_\beta} \cdot \max_s D_{TV}(\pi \parallel \hat{\pi}_\beta)[s] \cdot \mathbb{E}_{s \sim \rho_{\hat{\pi}_\beta}(\cdot)} [D_{TV}(\pi \parallel \hat{\pi}_\beta)[s]] \\ - 2\gamma \mathbb{A}_{\hat{\pi}_\beta} \cdot \max_s D_{TV}(\pi \parallel \hat{\pi}_\beta)[s] \cdot \mathbb{E}_{s \sim \rho_{\mathcal{D}}(\cdot)} [1 - \hat{\pi}_\beta(a | s)],$$



## □ BPPO - Method

### ➤ Stage 1: Supervised Learning



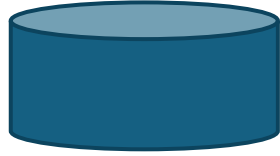
### ➤ Stage 2: Policy improvement using PPO

$$L_k(\pi) = \mathbb{E}_{s \sim \rho_{\mathcal{D}}(\cdot), a \sim \pi_k(\cdot|s)} \left[ \min \left( \frac{\pi(a|s)}{\pi_k(a|s)} \boxed{A_{\pi_k}(s, a)}, \text{clip} \left( \frac{\pi(a|s)}{\pi_k(a|s)}, 1 - 2\epsilon, 1 + 2\epsilon \right) A_{\pi_k}(s, a) \right) \right]$$

PPO objective with advantage replacement:  $A_\pi(s, a) \approx Q_{\hat{\pi}_\beta} - V_{\hat{\pi}_\beta}$

## □ RL partitioning

### ✓ Data interaction patterns



- On-policy RL
- off-policy RL
- Offline RL

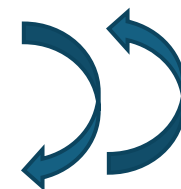
➤ A general version of RL policy update forms

Policy evaluation

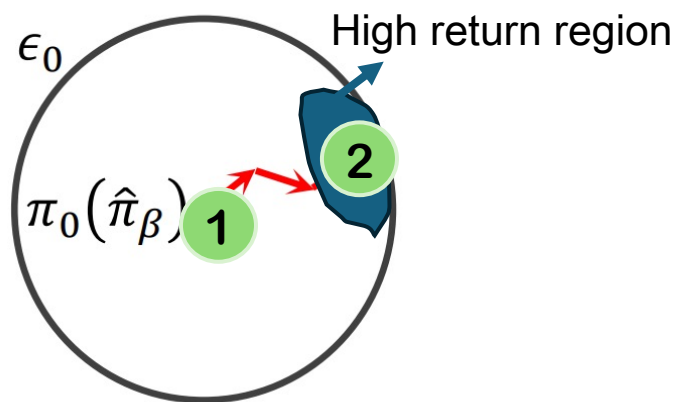
$$\hat{Q}^{k+1} = \arg \min_Q \mathbb{E}_{(s,a) \sim w} \left[ \left( \underbrace{T^{\pi^k} \hat{Q}}_{r + \gamma \mathbb{E}_{s'|s,a} \mathbb{E}_{a' \sim \pi^k(\cdot|s')} [\hat{Q}(s', a')]} - Q(s, a) \right)^2 \right]. \dots \textcircled{1}$$

Policy improvement

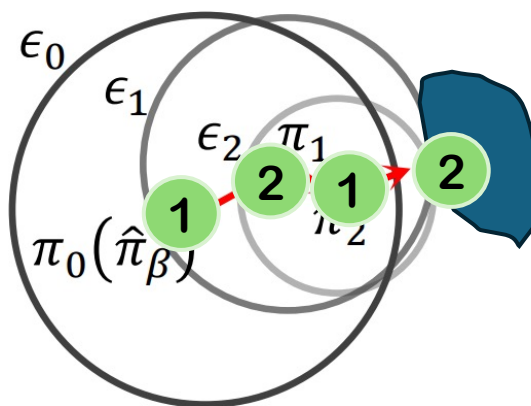
$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{s \sim \rho, a \sim \pi} [\hat{Q}^{k+1}(s, a)] - \beta \mathbb{E}_{s \sim \rho} [D_{\text{KL}}(\pi \parallel \pi_{\text{ref}})]. \dots \textcircled{2}$$



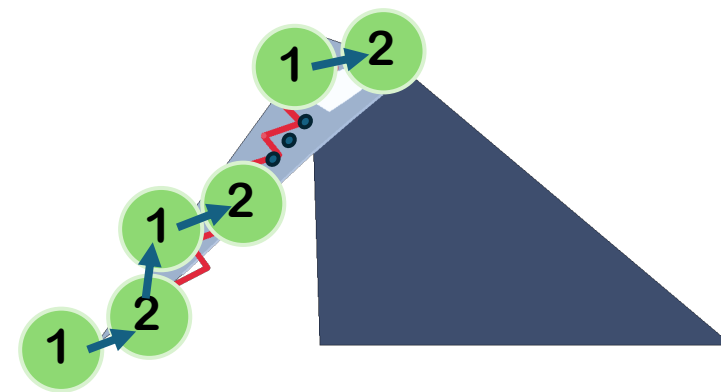
➤ One-step RL



➤ Multi-step RL



➤ Iterative RL



- One-step RL/IQL/IDQL...

- BPPO → Uni-O4/RL-100...

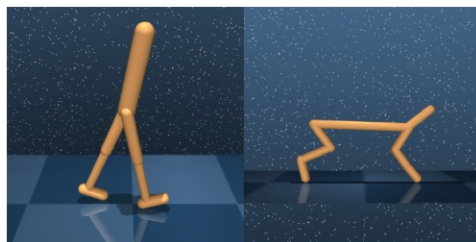
- PPO/SAC/...

More conservative, but stable

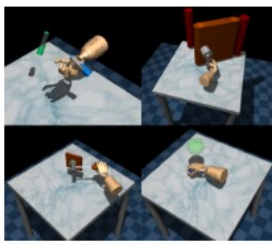
Aggressive exploration, but occasionally crash



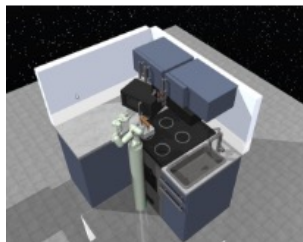
# Environments & Main Results



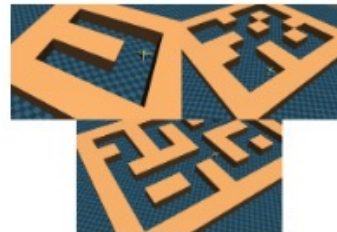
Gym locomotion



Adroit



Kitchen



Antmaze

Suite	Environment	Iterative methods		Onestep methods		BC (Ours)	BPPO (Ours)
		CQL	TD3+BC	Onestep RL	IQL		
Gym	halfcheetah-medium-v2	44.0	<b>48.3</b>	<b>48.4</b>	<b>47.4</b>	43.5±0.1	44.0±0.2
	hopper-medium-v2	58.5	59.3	59.6	66.3	61.3±3.2	<b>93.9±3.9</b>
	walker2d-medium-v2	72.5	<b>83.7</b>	81.8	78.3	74.2±4.6	<b>83.6±0.9</b>
	halfcheetah-medium-replay-v2	<b>45.5</b>	<b>44.6</b>	38.1	<b>44.2</b>	40.1±0.1	41.0±0.6
	hopper-medium-replay-v2	95.0	60.9	<b>97.5</b>	94.7	66.0±18.3	92.5±3.4
	walker2d-medium-replay-v2	77.2	<b>81.8</b>	49.5	73.9	33.4±11.2	77.6±7.8
	halfcheetah-medium-expert-v2	91.6	90.7	<b>93.4</b>	86.7	64.4±8.5	<b>92.5±1.9</b>
	hopper-medium-expert-v2	105.4	98.0	103.3	91.5	64.9±7.7	<b>112.8±1.7</b>
	walker2d-medium-expert-v2	108.8	110.1	<b>113.0</b>	109.6	107.7±3.5	<b>113.1±2.4</b>
	<i>Gym locomotion-v2 total</i>	698.5	677.4	684.6	692.4	555.5±57.2	<b>751.0±21.8</b>
Adroit	pen-human-v1	37.5	8.4*	90.7*	71.5	61.6±9.7	<b>117.8±11.9</b>
	hammer-human-v1	4.4	2.0*	0.2*	1.4	2.0±0.9	<b>14.9±3.2</b>
	door-human-v1	9.9	0.5*	-0.1*	4.3	7.8±3.5	<b>25.9±7.5</b>
	relocate-human-v1	0.2	-0.3*	2.1*	0.1	0.1±0.0	<b>4.8±2.2</b>
	pen-cloned-v1	39.2	41.5*	60.0	37.3	58.8±16.0	<b>110.8±6.3</b>
	hammer-cloned-v1	2.1	0.8*	2.0	2.1	0.5±0.2	<b>8.9±5.1</b>
	door-cloned-v1	0.4	-0.4*	0.4	1.6	0.9±0.8	<b>6.2±1.6</b>
	relocate-cloned-v1	-0.1	-0.3*	-0.1	-0.2	-0.1±0.0	<b>1.9±1.0</b>
	<i>adroit-v1 total</i>	93.6	52.2	155.2	118.1	131.6±31.1	<b>291.4±38.8</b>
Kitchen	kitchen-complete-v0	43.8	0.0*	2.0*	62.5	55.0±11.5	<b>91.5±8.9</b>
	kitchen-partial-v0	49.8	22.5*	35.5*	46.3	44.0±4.9	<b>57.0±2.4</b>
	kitchen-mixed-v0	51.0	25.0*	28.0*	51.0	45.0±1.6	<b>62.5±6.7</b>
	<i>kitchen-v0 total</i>	144.6	47.5	65.5	159.8	144.0±18.0	<b>211.0±18.0</b>
<i>locomotion+kitchen+adroit</i>		936.7	777.1	905.3	970.3	831.1±106.3	<b>1253.4±78.6</b>

↑ 50% performance  
based on BC

## ❑ How BPPO avoid overestimating:

- Once policy evaluation
- CLIP function in PPO loss

## ❑ Issues in BPPO

- The needs of online policy evaluation
- The performance is highly related to the estimated behavior policy
- Just offline



Uni-O4

# The post-training of robot (foundation) model

Offline reinforcement  
learning  
- Pretraining/Finetuning

## **Offline-to-online finetuning**

- for fast adaptation
- with safety consideration
- for task acceleration

Finetuning from multi-  
modal perception



# Uni-O4: Unifying **Online** and **Offline** Deep Reinforcement Learning with Multi-Step **On-Policy** Optimization

Kun Lei<sup>1</sup> Zhengmao He<sup>14</sup> Chenhao Lu<sup>2</sup> Kaizhe Hu<sup>12</sup> Yang Gao<sup>123</sup> Huazhe Xu<sup>123</sup>

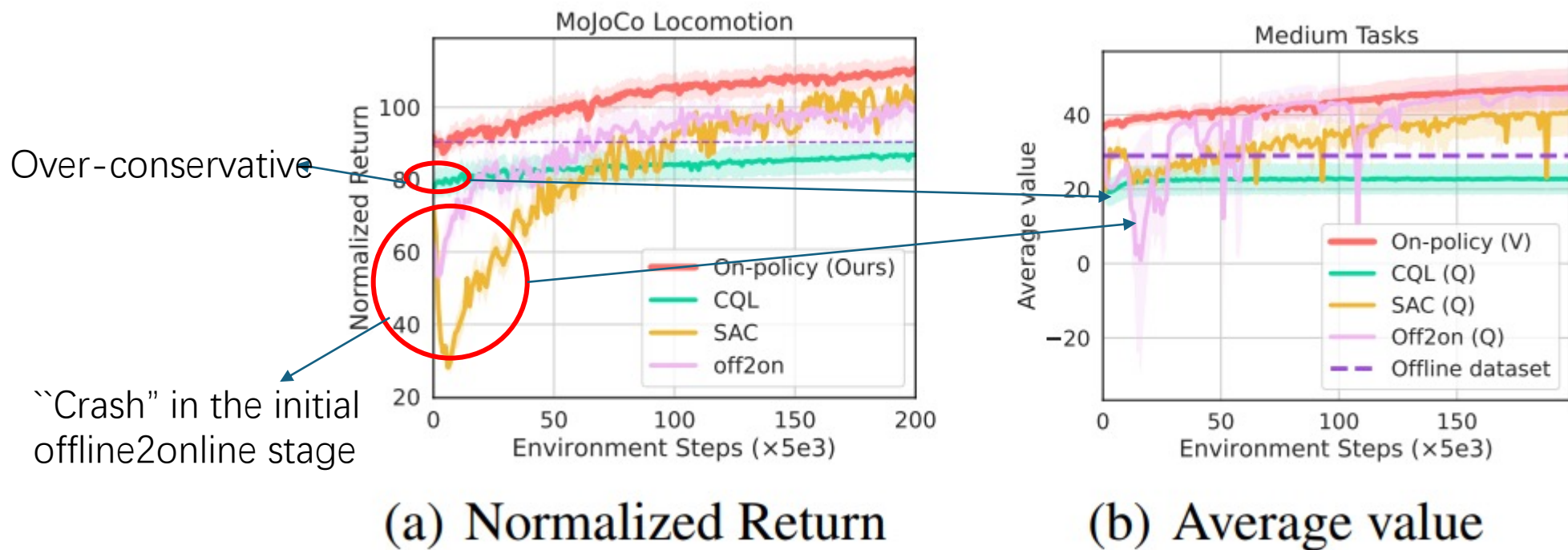
1 Shanghai Qi Zhi Institute. 2 Tsinghua University, IIS. 3 Shanghai AI Lab.

4 The Hong Kong University of Science and Technology (Guangzhou).



# Uni-O4: Unifying Online and Offline Deep Reinforcement Learning with Multi-Step On-Policy Optimization

**Offline-to-online issues:** *distribution shift* due to the *conservatism* used in offline phase.

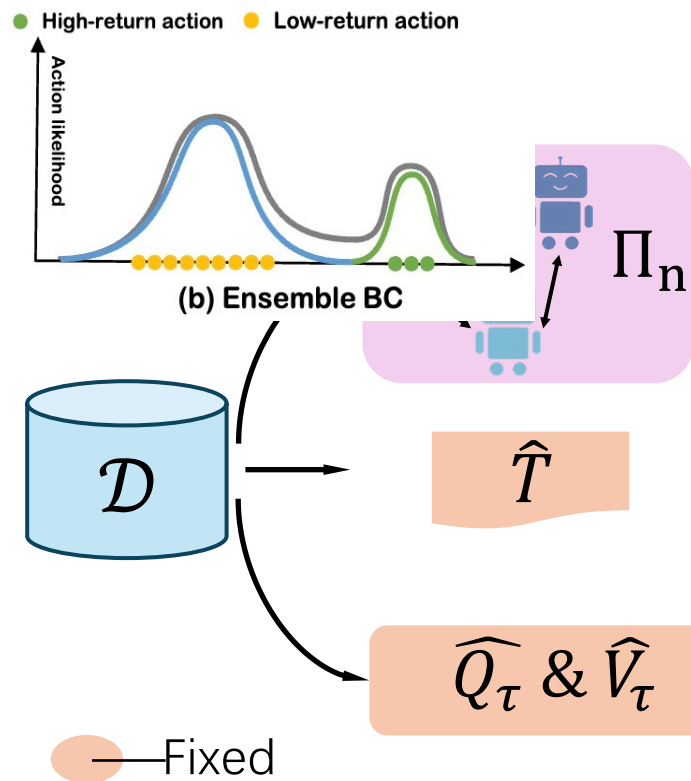


Can we **avoid introducing the conservatism** term during offline training and **eliminate the need for off-policy evaluation** during offline-to-online fine-tuning?

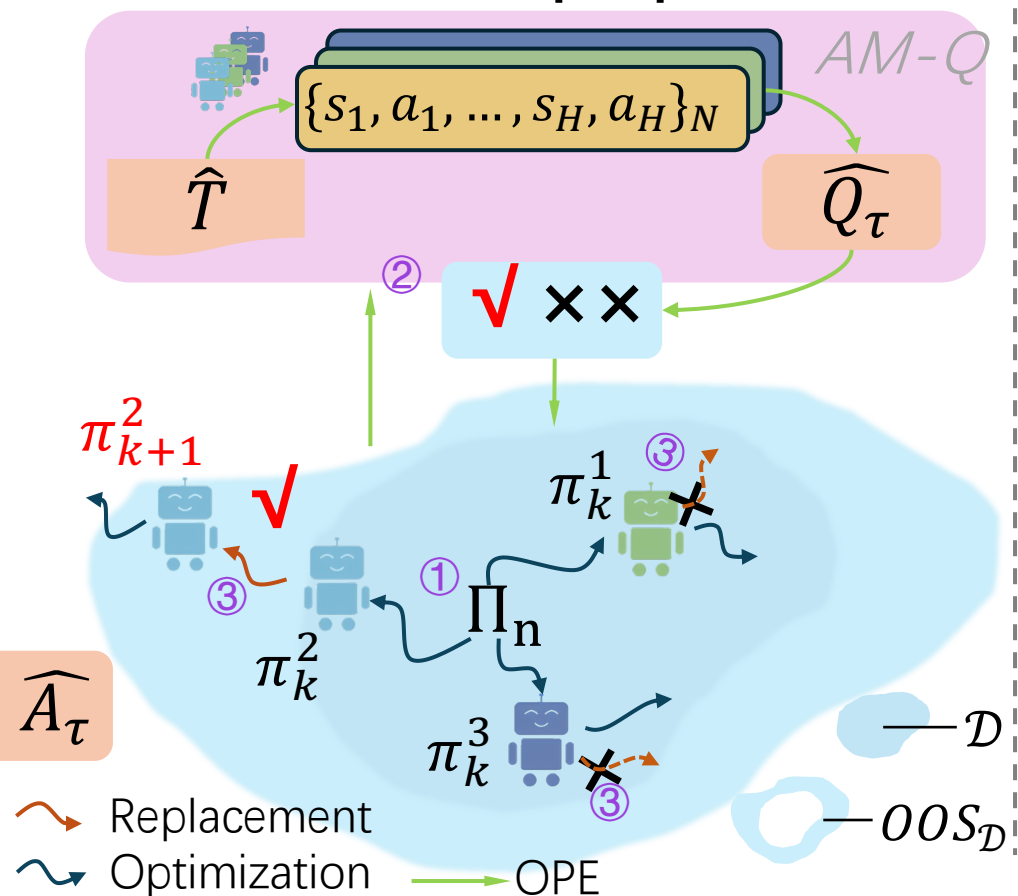
On-policy:  
Uni-O4 [1]

## ❑ Offline-to-online framework

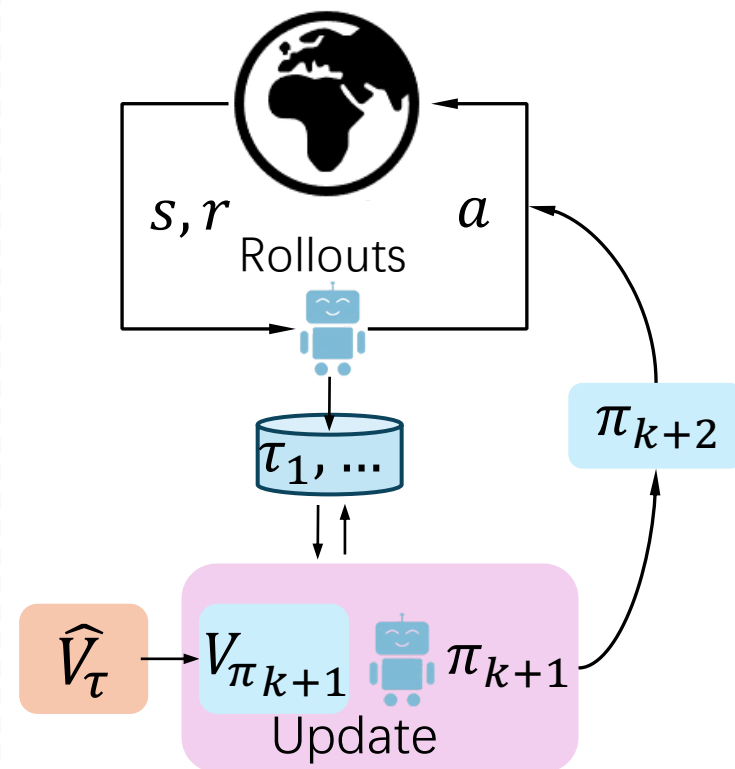
# Supervised Learning



## Offline Multi-Step Optimization

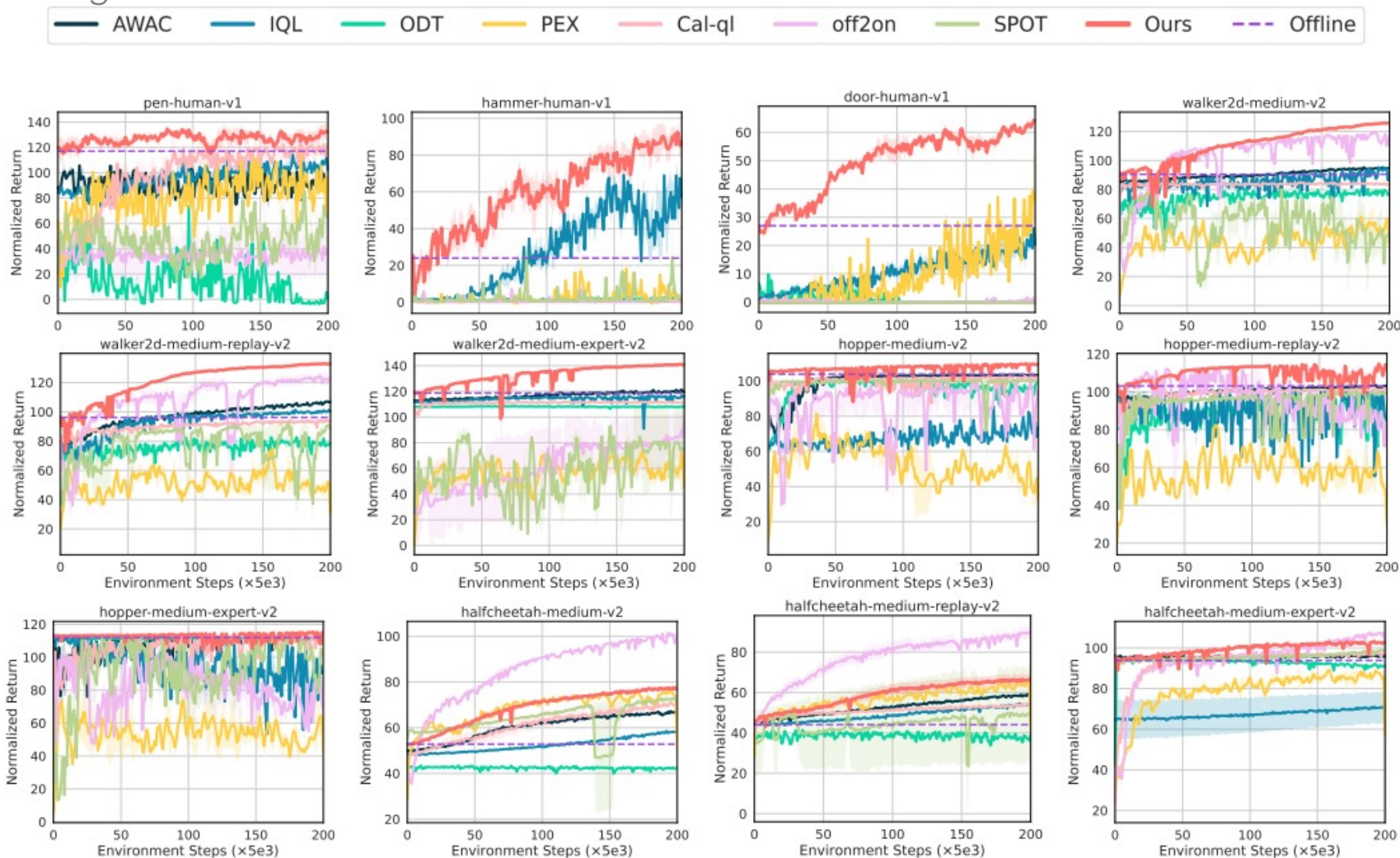


# Online Fine-Tuning

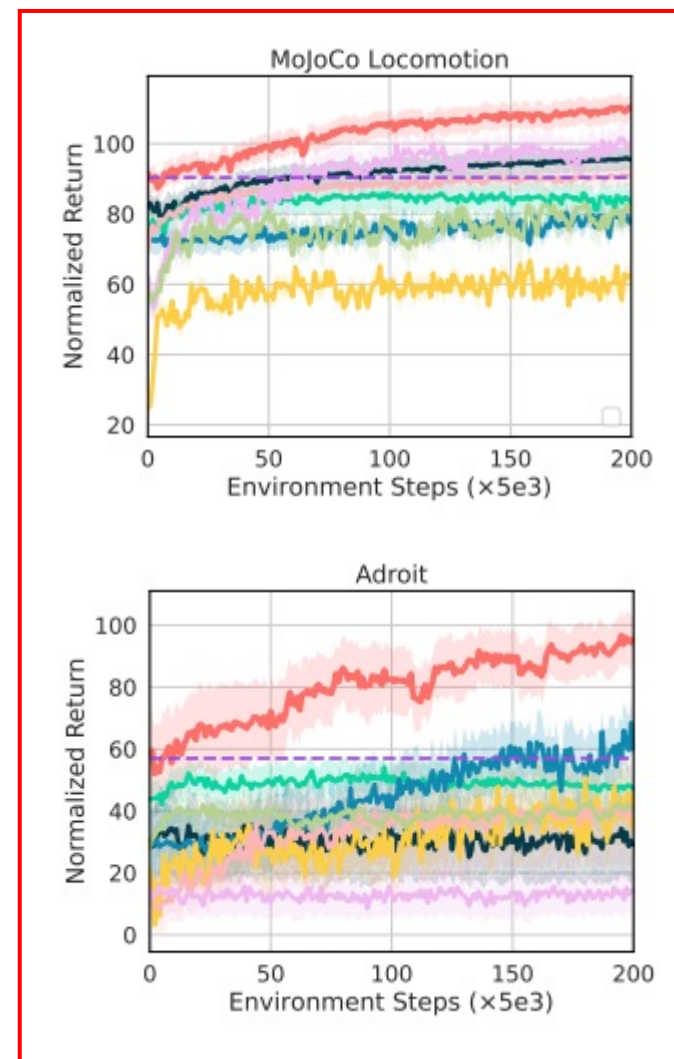


# □ Simulated Results

- Offline phase, offline-to-online phase, real robot setting



Results over single task

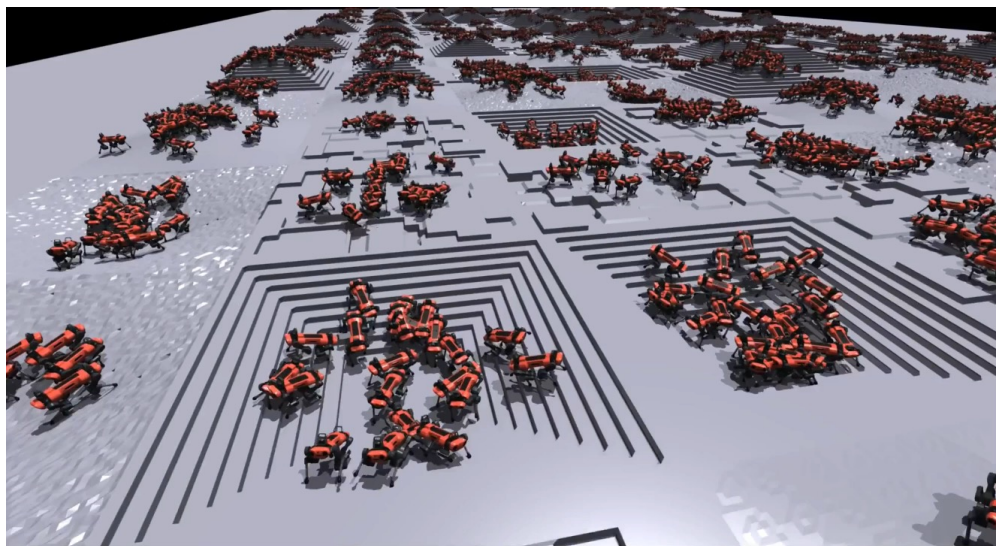


Average results over domains



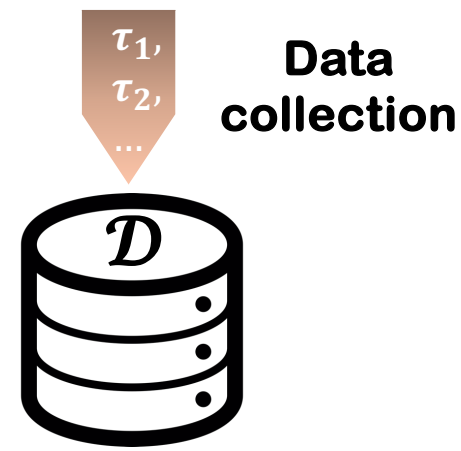
- Offline phase, offline-to-online phase, real robot setting

## Online Phase (Sim)



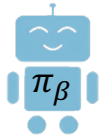
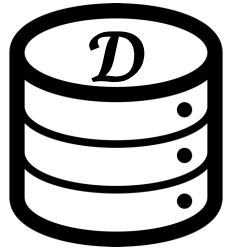
**Train quadruped robot  
For several minutes**

**Deploy in  
Real-World**





## Offline fine-tuning Phase (Real-world)



Supervised Learning

$\hat{T}$

$\widehat{Q}_\tau$  &  $\widehat{V}_\tau$

Offline Multi-Step Optimization

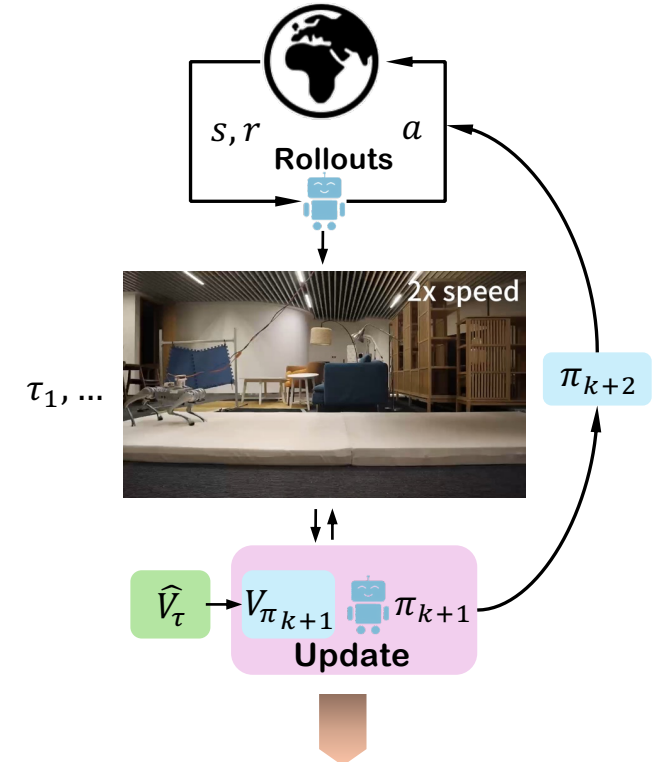


Performance Improved on  
soft & deformable terrain



But still **not satisfactory**  
when speed is **fast**

## Online Phase (Real-world)



**Faster** and more robust

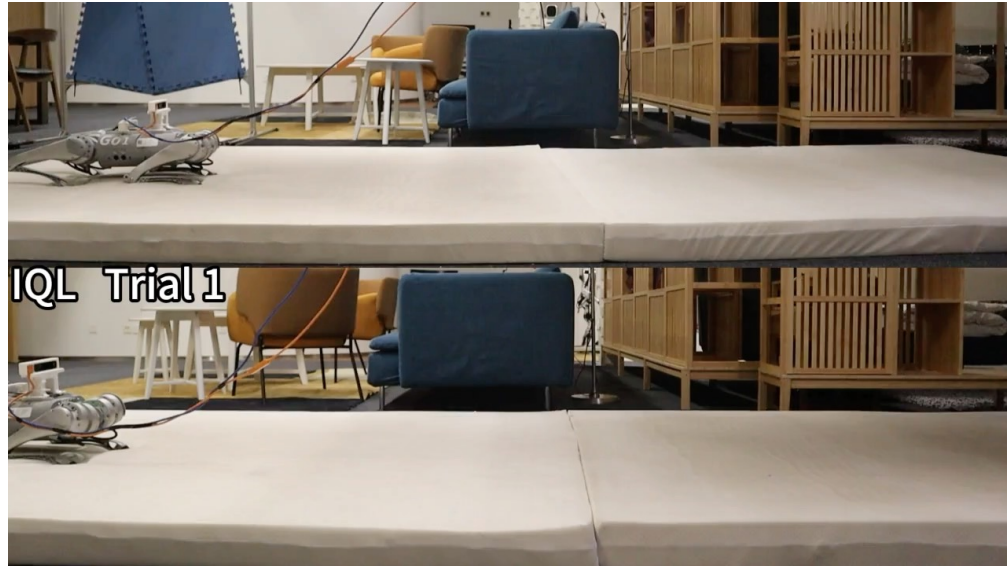
**Online Phase (Sim)**

**Offline Phase (Real-world)**

**Online Phase (Real-world)**

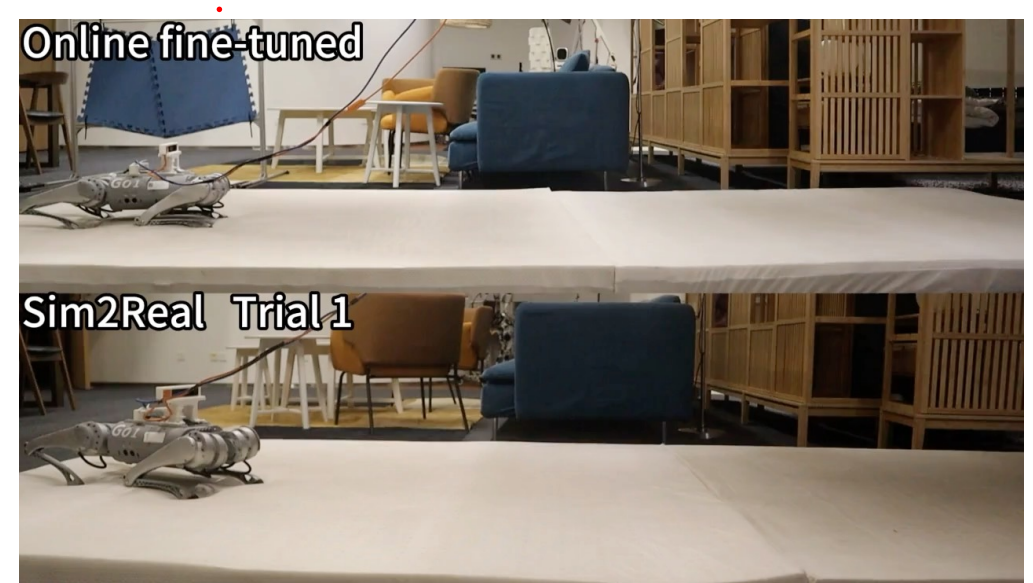
**All in one with PPO!**

Offline and online fine-tuning comparison



Uni-O4 vs. IQL [1]

Offline-to-online baseline



Uni-O4 vs. Walk these ways [2]

[1] Kostrikov I, Nair A, Levine S. Offline reinforcement learning with implicit q-learning[J]. arXiv preprint arXiv:2110.06169, 2021.

[2] Margolis, Gabriel B., and Pulkit Agrawal. "Walk these ways: Tuning robot control for generalization with multiplicity of behavior." *Conference on Robot Learning*. PMLR, 2023.

Key insight of Uni-O4:

- On-policy RL can unify offline and online setting
- Offline RL could work as a finetuning paradigm

Only explored locomotion tasks.

# The post-training of robot (foundation) model

Offline reinforcement  
learning  
- Pretraining/Finetuning

Offline-to-online finetuning  
- for fast adaptation  
- with safety consideration  
- for task acceleration

**Finetuning from  
multi-modal perception**



## RL-100: Performant Robotic Manipulation with Real-World Reinforcement Learning

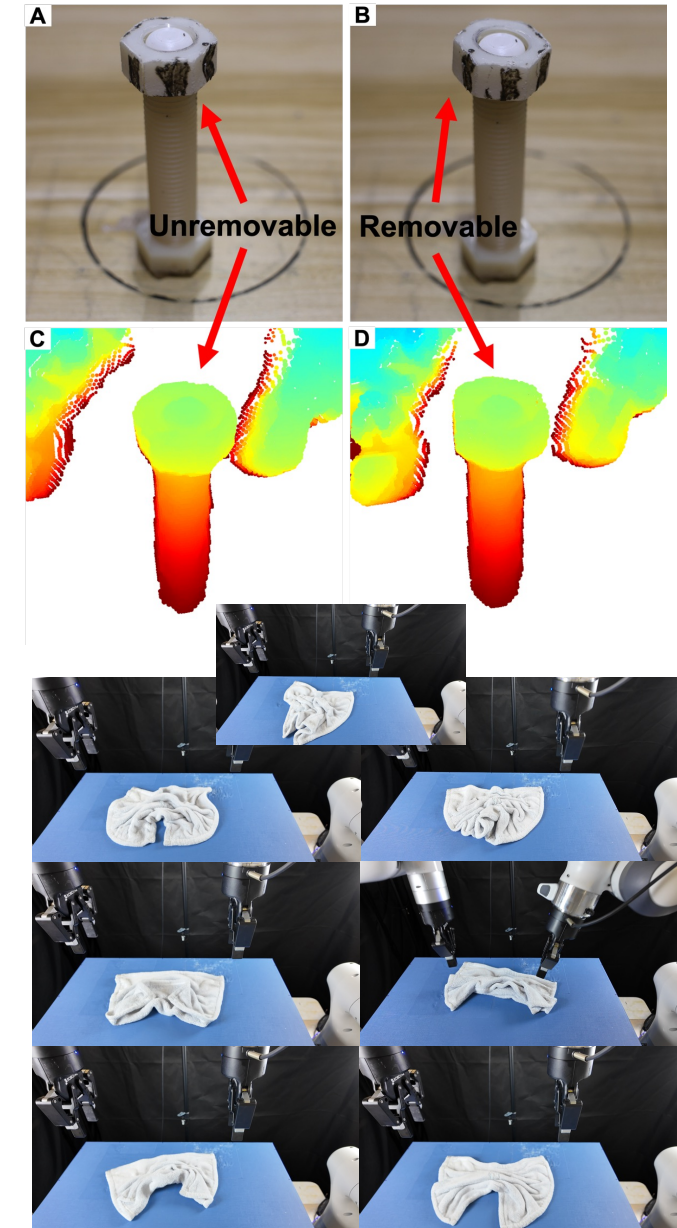
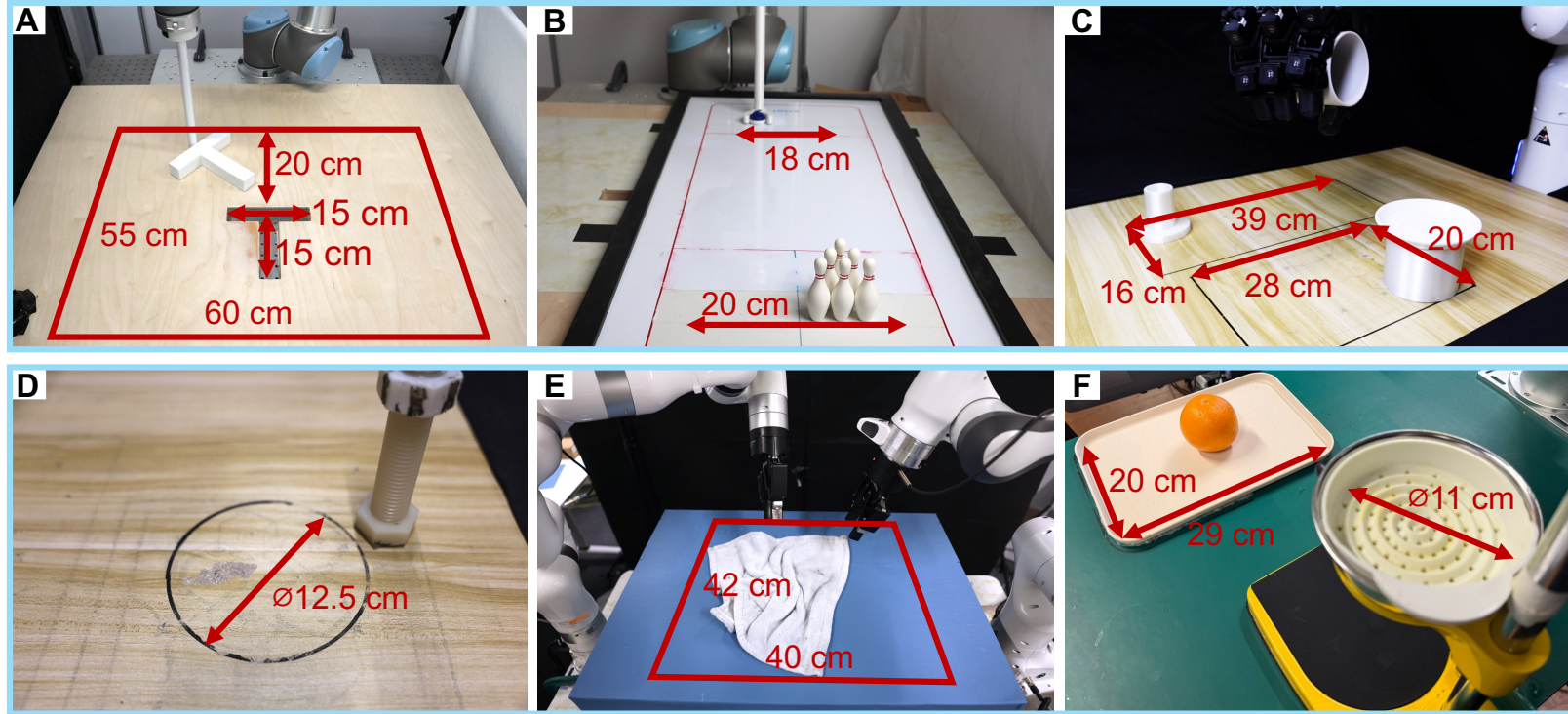
**Kun Lei\***, Huanyu Li\*, Dongjie Yu\*, Zhenyu Wei\*, Lingxiao Guo, Zhennan Jiang, Ziyu Wang, Shiyu Liang, Huazhe Xu.

Preparing to submit.



**What we did?**

- **7** real robot tasks, **900/900** successes. Up to **250** consecutive trials in one task, running **2 hours** nonstop without failure.
- High success rate against physical disturbances, zero-shot, and few-shot adaptation



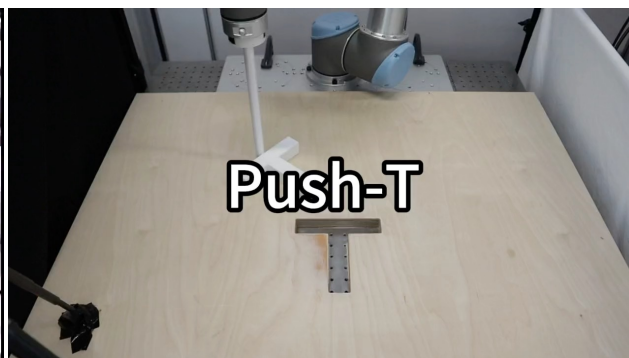


What we did?



Folding

- Folding
  - Dual-arm
  - Deformable



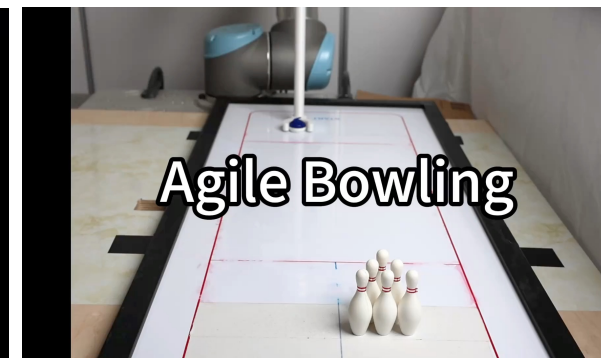
Push-T

- Dynamic Push-T
  - Rigid-body dynamics



Nuts Pouring

- Pour
  - Fluids / granular



Agile Bowling

- Bowling
  - agile



Orange-juicing  
Placing

- Juicing-stage 1
  - Diff. Size
  - Various inclination



Orange-juicing  
Removal

- Juicing-stage 2
  - Deformable
  - Confined-space



Dynamic  
Unscrewing

- Unscrewing
  - Dynamic



- Serve for 7 hours
  - Zero-shot

■ Key Words: ➤ 250/250 ➤ 2 Hours ➤ Efficient ➤ 7 hours outdoor serving

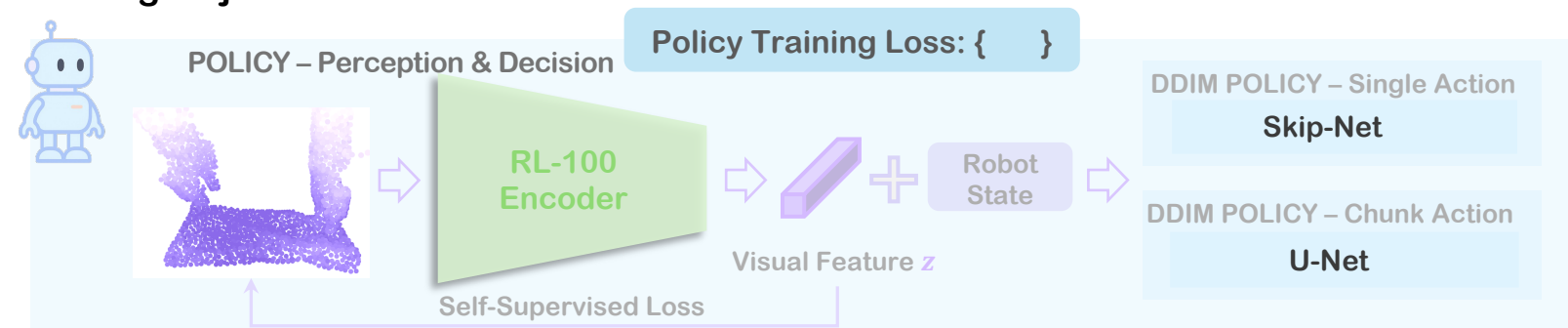
**How to do it**



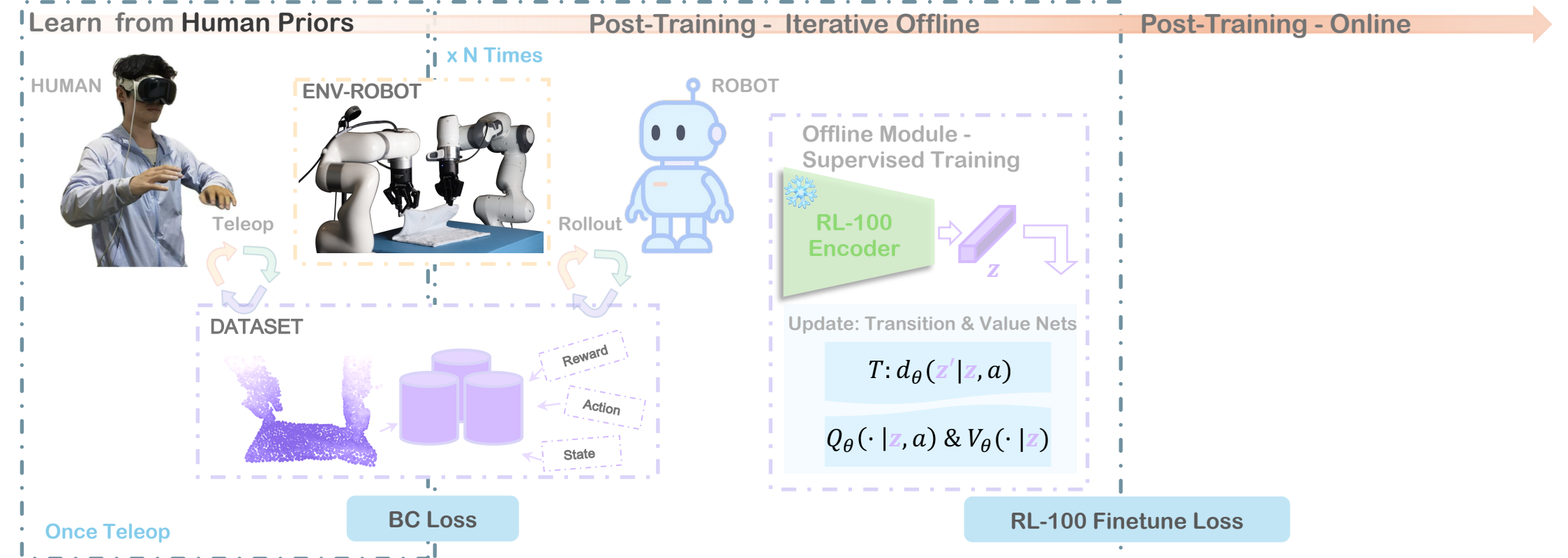
## Training Pipeline



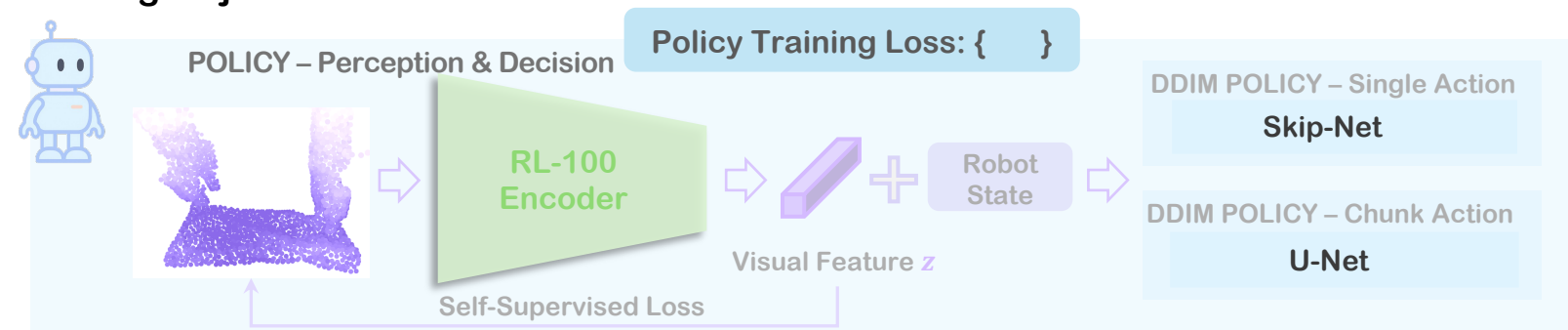
## Training Objective



## Training Pipeline



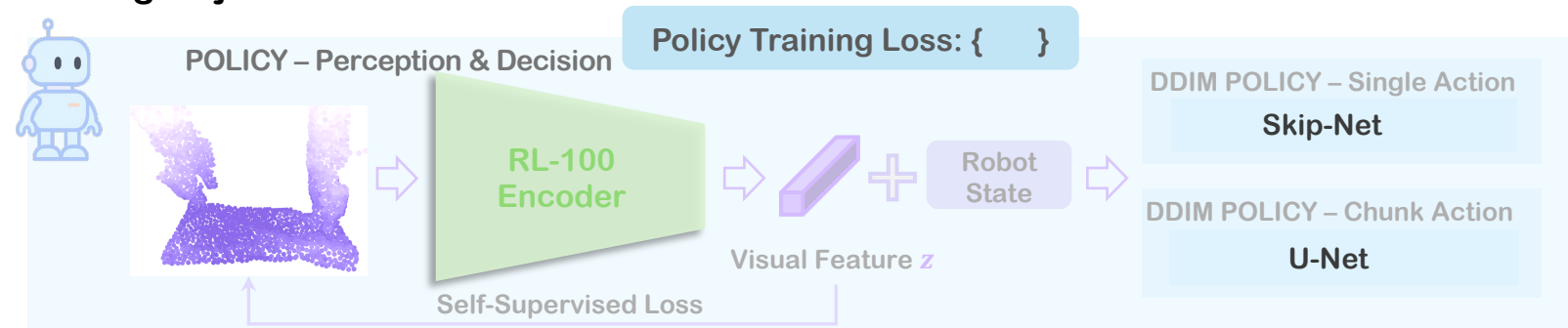
## Training Objective



## Training Pipeline



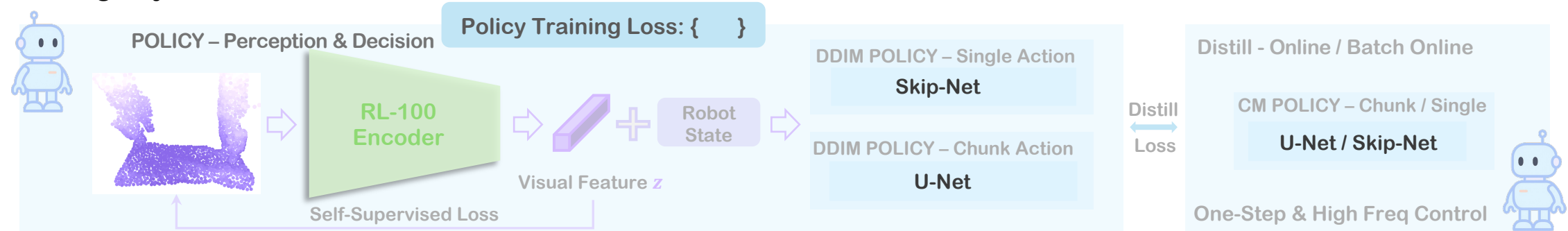
## Training Objective



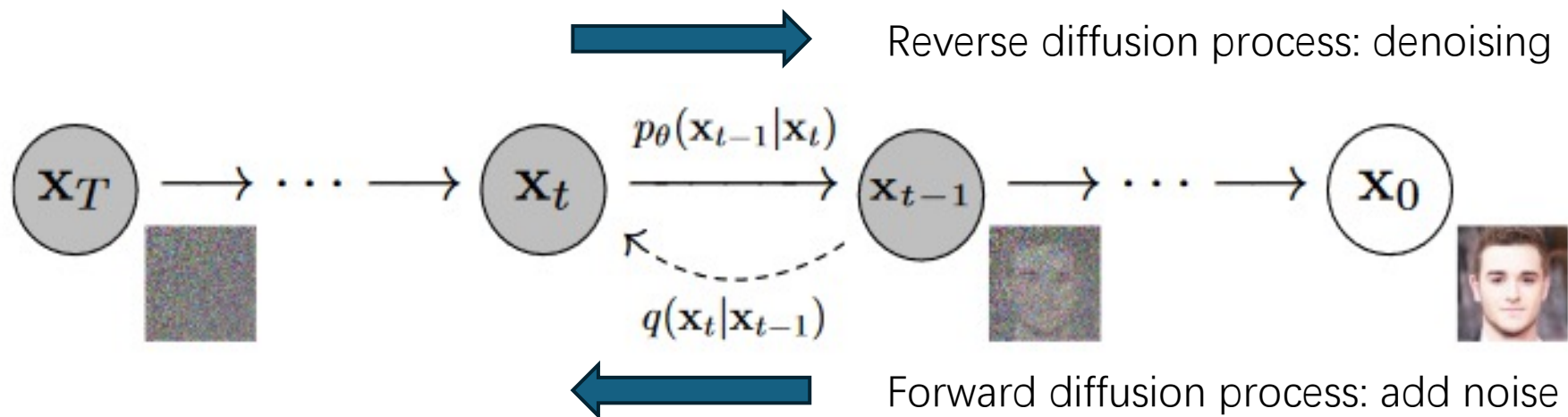
## Training Pipeline



## Training Objective



▣ Two-level MDP with DDIM sampling with stochastic form



Policy gradient (PG) loss:  $\nabla_\theta J = \mathbb{E}_\pi[\nabla_\theta \log p_\theta(a|s) R_\pi]$

PG loss with multi step sampling (DDIM):  $\nabla_\theta J = \mathbb{E}_\pi\left[\sum_{k=0}^K \nabla_\theta \log p_\theta(x_{\tau_{k-1}}|x_{\tau_k}, s) R_\pi\right]$

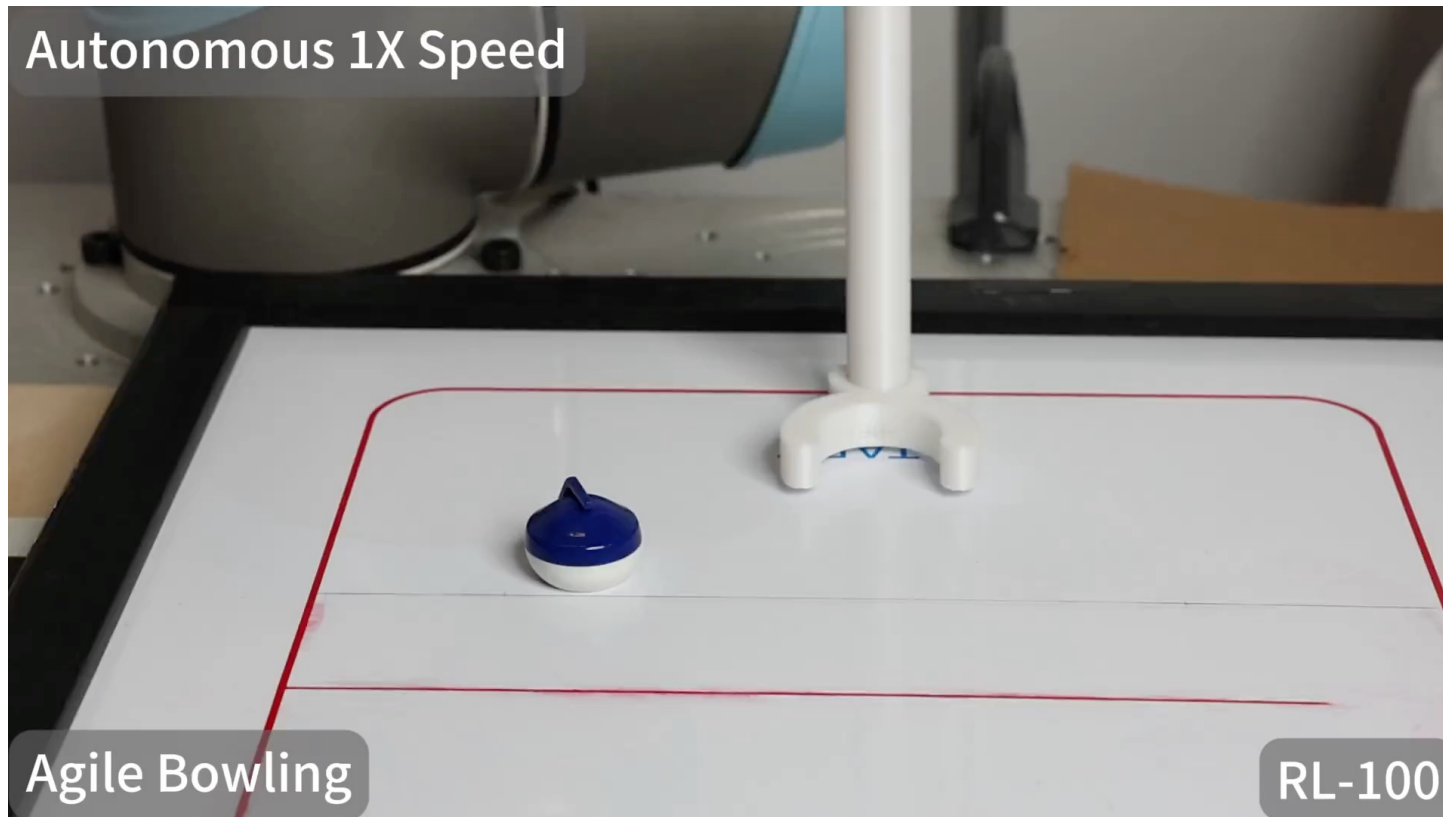
PG loss with multi step sampling and importance sampling:  $\nabla_\theta J = \mathbb{E}_\pi\left[\sum_{k=0}^K \nabla_\theta \log \frac{p_\theta(x_{\tau_{k-1}}|x_{\tau_k}, s)}{p_{\theta_{old}}(x_{\tau_{k-1}}|x_{\tau_k}, s)} A_\pi\right]$   $\left\{ \begin{array}{l} A_\pi = Q - V \\ \text{offline} \\ A_\pi = GAE \\ \text{online} \end{array} \right.$

One-step consistency distillation:  $\mathcal{L}_{CD}(\theta) = \mathbb{E}_{x_0, \tau, \varepsilon} \left[ \left\| C_\theta(x^\tau, \tau) - \text{sg}[\Psi_\varphi(x^\tau, \tau \rightarrow 0)] \right\|_2^2 \right]$

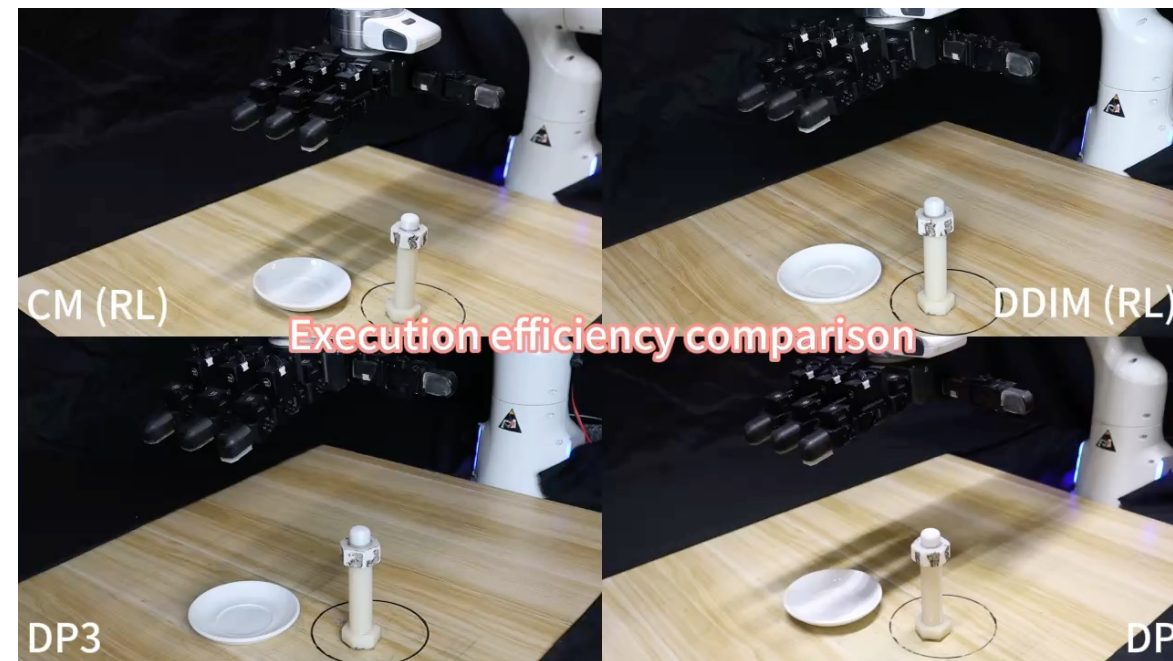
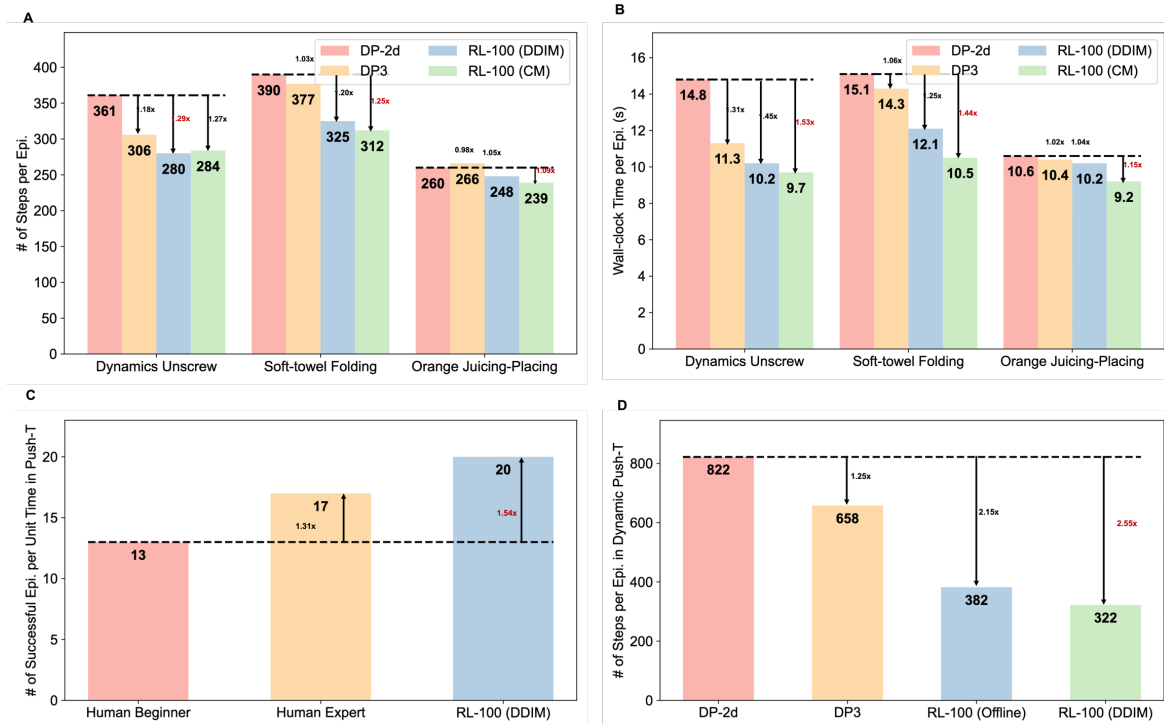
Overall finetune loss:  $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{RL}} + \lambda_{\text{CD}} \cdot \mathcal{L}_{\text{CD}}$



- 
- Robustness, Zero-Shot & Few-shot Generalization



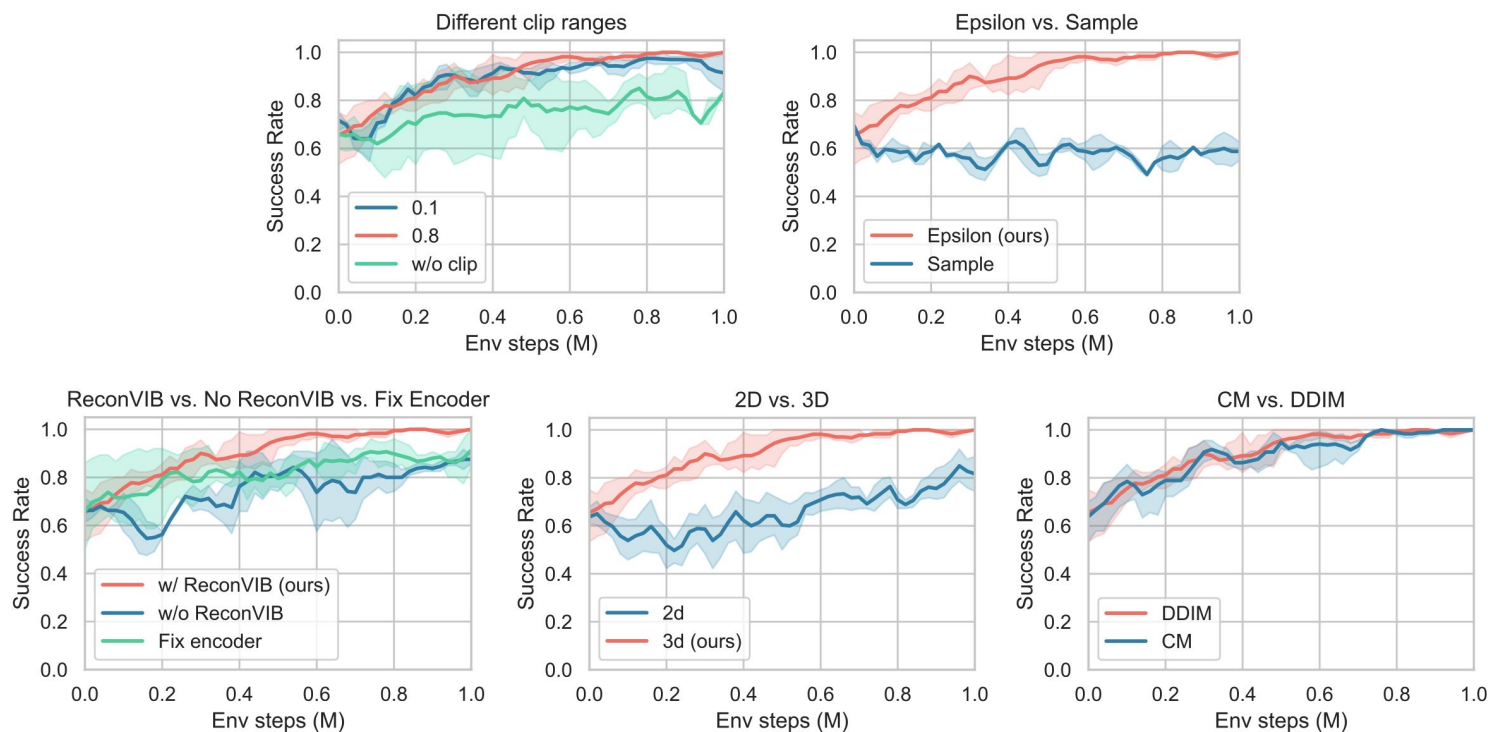
## Execution efficiency



□ Data usage

Task	Human Demonstration		Iterative Offline RL		Online RL	
	# of epi.	Collection time (h)	# of epi.	Collection time (h)	# of epi.	Collection time (h)
Dynamic Push-T	100	2	821	8	763	7.5
Agile Bowling	100	2	249	2	213	2.5
Pouring	64	1	741	6.8	129	1.5
Soft-towel Folding	400	5	896	11	654	8.5
Dynamic Unscrew	31	0.5	467	4.5	288	3
Orange Juicing – Placing	80	1.5	642	10.5	750	12.5
Orange Juicing – Removal	29	0.5	149	2.5	240	4
<b>Average</b>	<b>115</b>	<b>1.8</b>	<b>566</b>	<b>6.5</b>	<b>434</b>	<b>5.6</b>

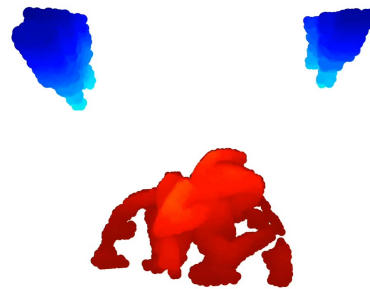
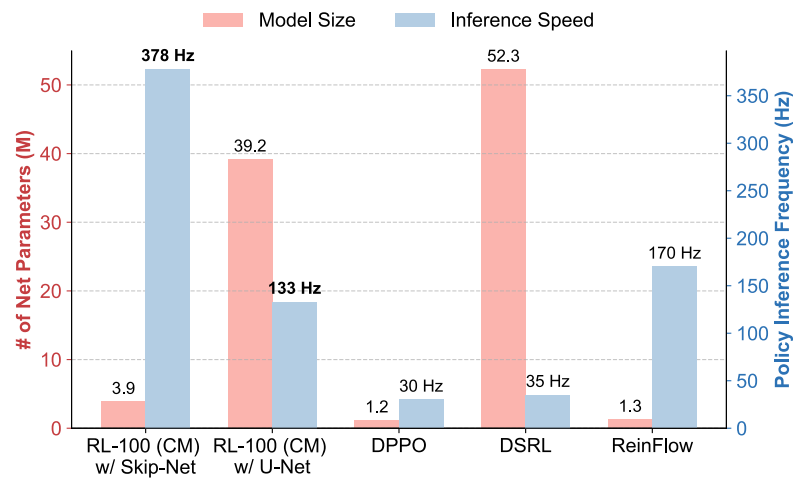
## □ Ablation study



Takeaway:

- 1) **Variance clipping** is valid for **stable exploration** - variance clipping in the stochastic DDIM sampling process.
- 2) **Epsilon prediction** is more suitable for RL: large noise schedule for exploration
- 3) **Reconstruction** is crucial for **visual robotic manipulation RL** as it mitigates representational drift and improves sample efficiency.
- 4) On a relatively clean scene, the 3D variant learns **faster and attains a higher** final success rate.
- 5) **CM** effectively compresses the iterative denoising process without sacrificing control quality, enabling **high-frequency deployment**.

# Ablation study



Folding





□ Next move - Liberate productive forces: robot helps

Single task



Multi task:  
The same series  
More data-more  
robots



Understand  
humans' instructions





# Thanks!



Project page



Wechat