

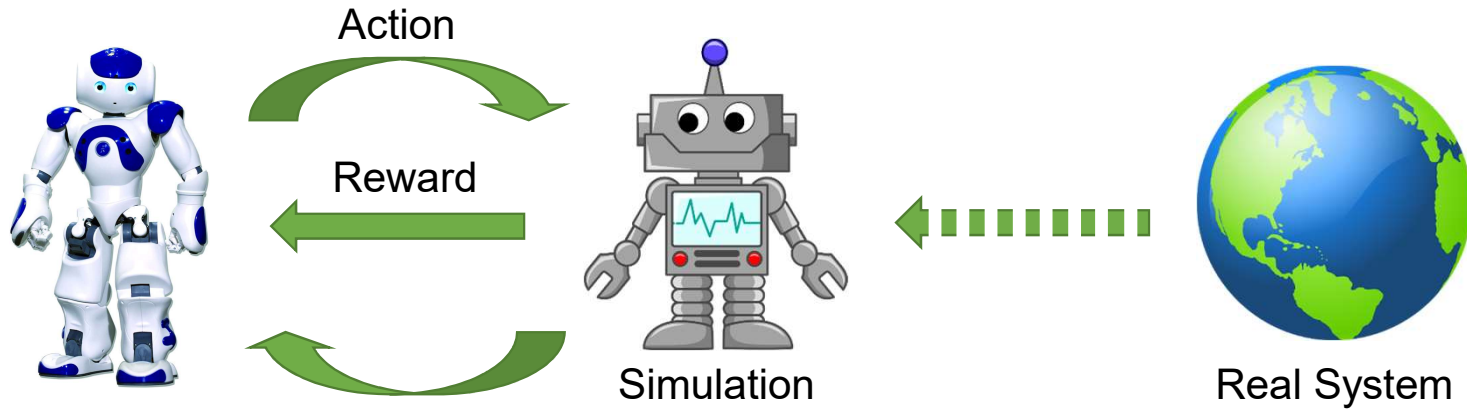


# When to Trust Your Simulator: Dynamics-Aware Hybrid Offline-and-Online Reinforcement Learning

Haoyi Niu, Shubham Sharma, Yiwen Qiu, Ming Li, Guyue Zhou, Jianming Hu, Xianyuan Zhan



# Real-World Challenges for RL



## Visual Gap

## Dynamics Gap

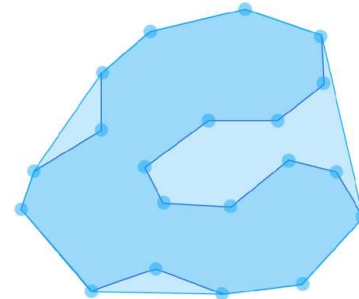
*Simulation*

*Real World*



*Convex Surface*

*Rigid Body*

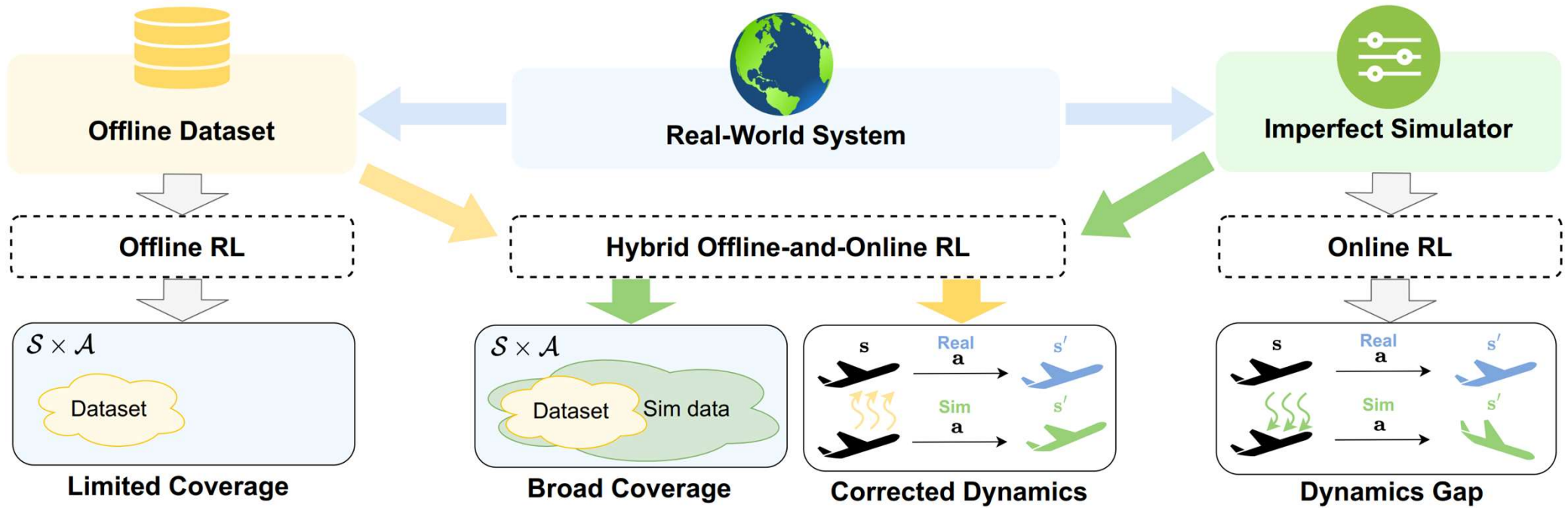


Domain Randomization  
Robust and transferable state-encoder  
proper model setups, etc.

causes a systematic bias & widely exist

**Not easy to address**

# Hybrid Offline-and-Online RL (H2O)



# Hybrid Offline-and-Online RL (H2O)

## Dynamics-Aware Policy Evaluation:

$$\min_Q \beta \left( \underbrace{\log \sum_{\mathbf{s}, \mathbf{a}} \omega(\mathbf{s}, \mathbf{a}) \exp(Q(\mathbf{s}, \mathbf{a}))}_{\text{Minimize the dynamics-gap weighted soft-maximum of Q values:}} - \underbrace{\mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} [Q(\mathbf{s}, \mathbf{a})]}_{\text{Maximize Q values on data:}} \right) + \tilde{\mathcal{E}}(Q, \hat{B}^\pi \hat{Q})$$

**Minimize the dynamics-gap weighted soft-maximum of Q values:**

Push down Q values on high dynamics-gap samples

**Maximize Q values on data:**

Pull up Q values on real offline data samples

$$\omega(\mathbf{s}, \mathbf{a}) = u(\mathbf{s}, \mathbf{a}) / \sum_{\tilde{\mathbf{s}}, \tilde{\mathbf{a}}} u(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})$$

$$u(\mathbf{s}, \mathbf{a}) := D_{KL}(P_{\hat{\mathcal{M}}} \| P_{\mathcal{M}})$$

$$\approx \sum_{s'_i \sim P_{\hat{\mathcal{M}}}(s'_i | \mathbf{s}, \mathbf{a})}^N \log \frac{P_{\hat{\mathcal{M}}}(s'_i | \mathbf{s}, \mathbf{a})}{P_{\mathcal{M}}(s'_i | \mathbf{s}, \mathbf{a})}$$

$$= \sum_{s'_i \sim P_{\hat{\mathcal{M}}}(s'_i | \mathbf{s}, \mathbf{a})}^N \log \left[ \frac{1 - p(\text{real} | \mathbf{s}, \mathbf{a}, s')}{p(\text{real} | \mathbf{s}, \mathbf{a}, s')} / \frac{1 - p(\text{real} | \mathbf{s}, \mathbf{a})}{p(\text{real} | \mathbf{s}, \mathbf{a})} \right]$$

$$= \sum_{s'_i \sim P_{\hat{\mathcal{M}}}(s'_i | \mathbf{s}, \mathbf{a})}^N \log \left[ \frac{1 - D_{sas}(\mathbf{s}, \mathbf{a}, s')}{D_{sas}(\mathbf{s}, \mathbf{a}, s')} / \frac{1 - D_{sa}(\mathbf{s}, \mathbf{a})}{D_{sa}(\mathbf{s}, \mathbf{a})} \right]$$

$$\tilde{\mathcal{E}}(Q, \hat{B}^\pi \hat{Q}) = \frac{1}{2} \mathbb{E}_{\mathbf{s}, \mathbf{a}, s' \sim \mathcal{D}} \left[ (Q - \hat{B}^\pi \hat{Q})(\mathbf{s}, \mathbf{a}) \right]^2 + \frac{1}{2} \mathbb{E}_{\mathbf{s}, \mathbf{a}, s' \sim \mathcal{B}} \left[ \frac{P_{\mathcal{M}}(s' | \mathbf{s}, \mathbf{a})}{P_{\hat{\mathcal{M}}}(s' | \mathbf{s}, \mathbf{a})} (Q - \hat{B}^\pi \hat{Q})(\mathbf{s}, \mathbf{a}) \right]^2$$

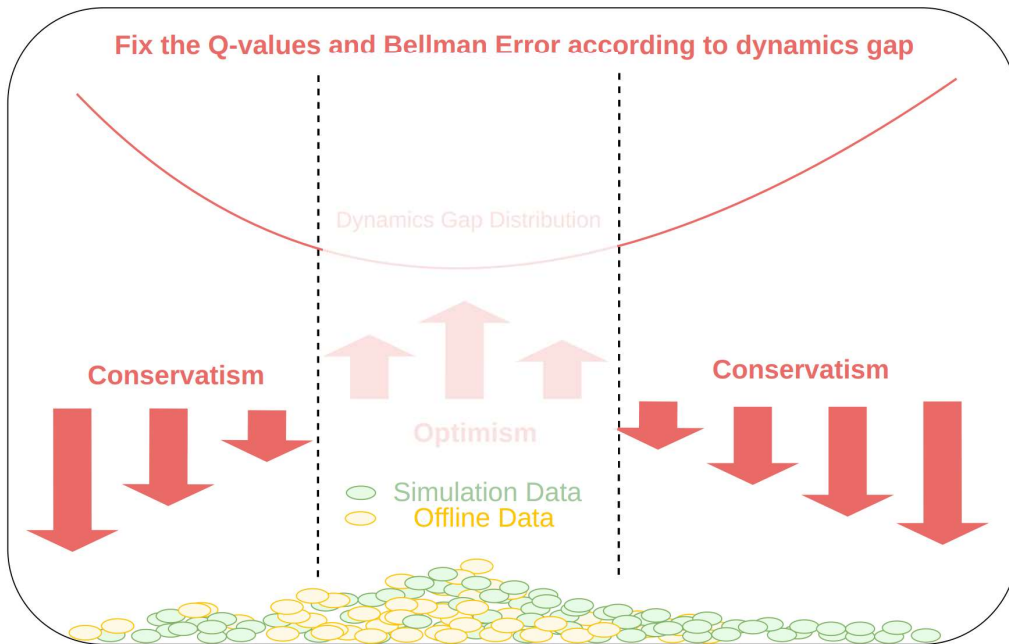
**Learn on both offline data and online simulated samples**

**Fix Bellman error due to dynamics gap:**

Use dynamics ratio as an importance sampling weight

# Theoretical Interpretation

$$\min_Q \beta \left( \log \sum_{\mathbf{s}, \mathbf{a}} \omega(\mathbf{s}, \mathbf{a}) \exp(Q(\mathbf{s}, \mathbf{a})) - \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} [Q(\mathbf{s}, \mathbf{a})] \right) + \tilde{\mathcal{E}}(Q, \hat{\mathcal{B}}^\pi \hat{Q})$$



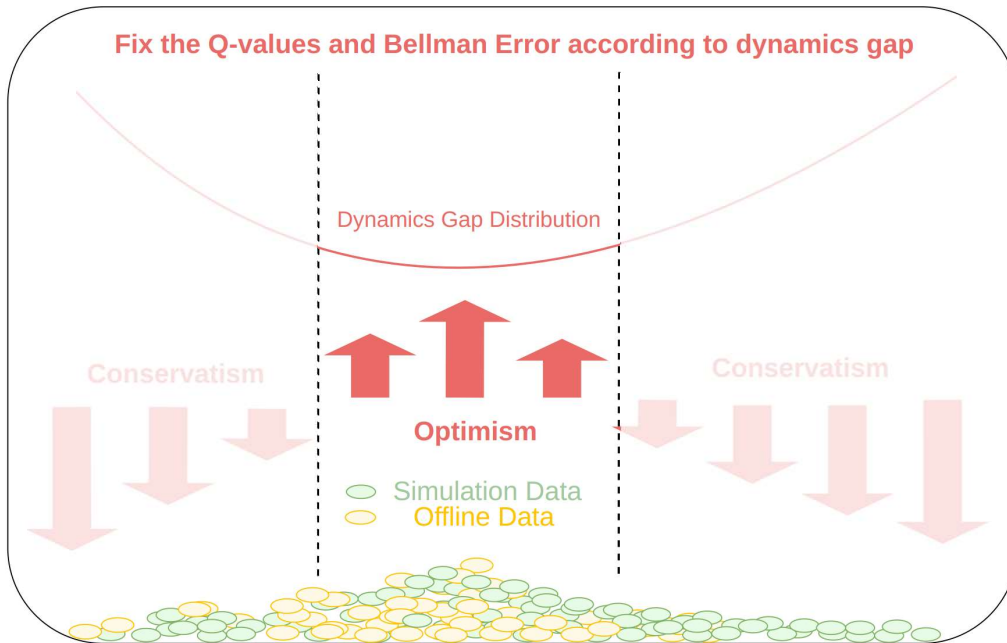
## Adaptation but not Conservatism

Can be interpreted as adding an adaptive adjustment on Q-values:

$$\hat{Q}^{k+1}(\mathbf{s}, \mathbf{a}) = (\hat{\mathcal{B}}^\pi \hat{Q}^k)(\mathbf{s}, \mathbf{a}) - \beta \left[ \frac{\omega(\mathbf{s}, \mathbf{a}) - d_M^{\pi_{\mathcal{D}}}(\mathbf{s}, \mathbf{a})}{d_M^{\pi_{\mathcal{D}}}(\mathbf{s}, \mathbf{a}) + d_M^{\pi_{\hat{\mathcal{D}}}}(\mathbf{s}, \mathbf{a})} \right]$$

# Theoretical Interpretation

$$\min_Q \beta \left( \log \sum_{\mathbf{s}, \mathbf{a}} \omega(\mathbf{s}, \mathbf{a}) \exp(Q(\mathbf{s}, \mathbf{a})) - \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} [Q(\mathbf{s}, \mathbf{a})] \right) + \tilde{\mathcal{E}}(Q, \hat{\mathcal{B}}^\pi \hat{Q})$$

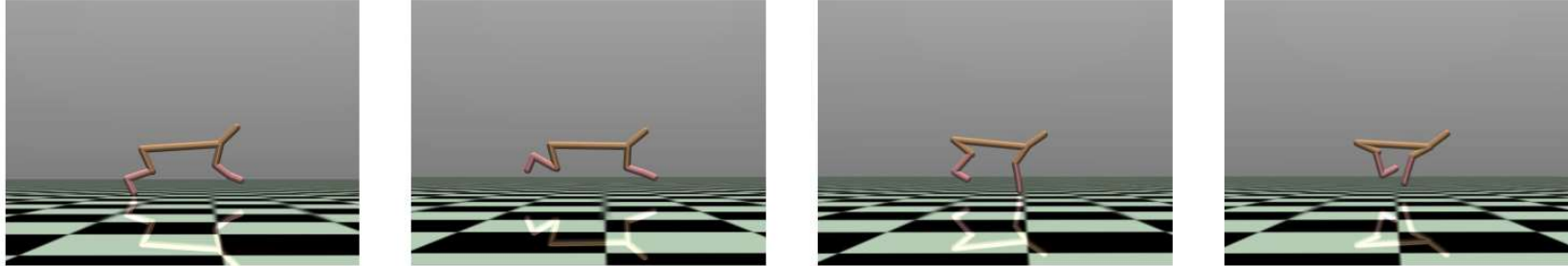


## Adaptation but not Conservatism

Can be interpreted as adding an adaptive adjustment on Q-values:

$$\hat{Q}^{k+1}(\mathbf{s}, \mathbf{a}) = (\hat{\mathcal{B}}^\pi \hat{Q}^k)(\mathbf{s}, \mathbf{a}) - \beta \left[ \frac{\omega(\mathbf{s}, \mathbf{a}) - d_M^{\pi_{\mathcal{D}}}(\mathbf{s}, \mathbf{a})}{d_M^{\pi_{\mathcal{D}}}(\mathbf{s}, \mathbf{a}) + d_M^{\pi_{\hat{\mathcal{D}}}}(\mathbf{s}, \mathbf{a})} \right]$$

# Simulation Experiments



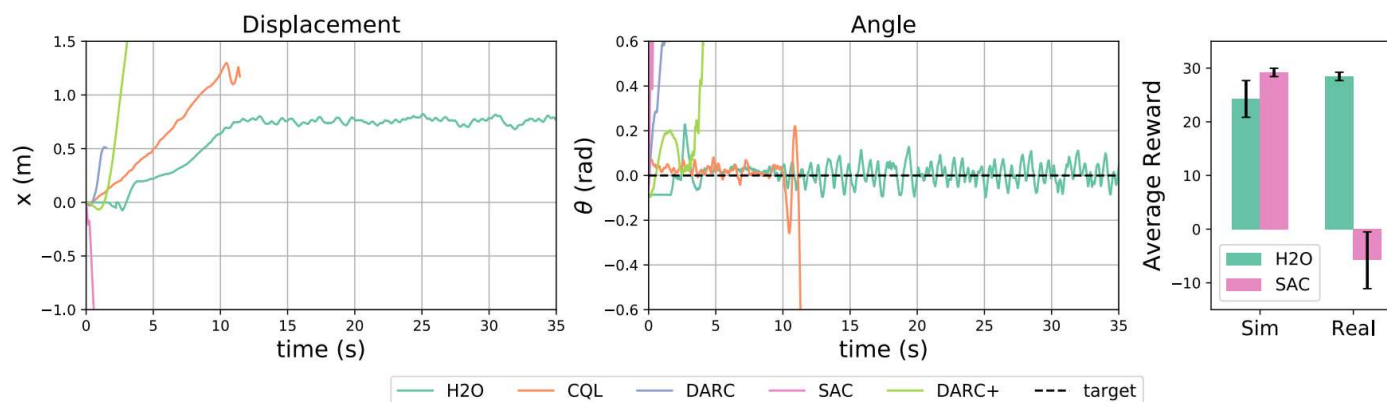
Dataset	Unreal Dynamics	SAC	CQL	DARC	DARC+	H2O
Medium	Gravity	4513±513	6066±73	5011±456	5706±440	<b>7085±416</b>
	Friction	2684±2646	6066±73	6113±104	6047±112	<b>6848±445</b>
	Joint Noise	4137±805	6066±73	5484±171	5314±520	<b>7212±236</b>
Medium Replay	Gravity	4513±513	5774±214	5105±460	4958±540	<b>6813±289</b>
	Friction	2684±2646	5774±214	5503±263	5288±100	<b>5928±896</b>
	Joint Noise	4137±805	5774±214	5137±225	5230±209	<b>6747±427</b>
Medium Expert	Gravity	4513±513	3748±892	<b>4759±353</b>	72±109	<b>4707±779</b>
	Friction	2684±2646	3748±892	<b>9038±1480</b>	7989±3999	6745±562
	Joint Noise	4137±805	3748±892	<b>5288±104</b>	733±767	<b>5280±1329</b>

## MuJoCO HalfCheetah Environment

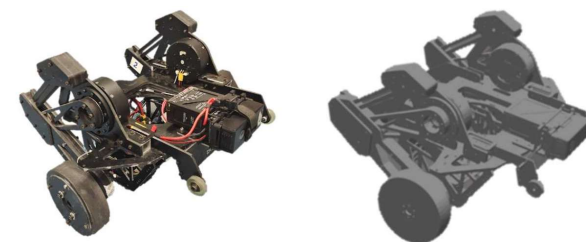
- **Offline dataset:** D4RL
- **Trained in simulation with unreal dynamics:** Gravity x2, Friction x0.3, Joint Noise N(0,1)
- **Evaluated in simulation with original dynamics:** Gravity x1, Friction x1, Joint Noise: none
- **Baseline:** SAC (online sim), CQL (offline real), DARC (online cross-domain, DARC+ (online+offline))

# Real-World Validation

## Standing Still



## Wheel-Legged Robot



SAC

DARC

DARC+

CQL

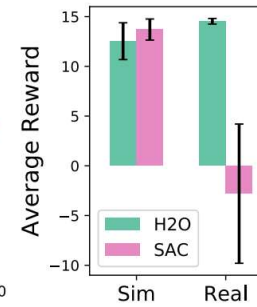
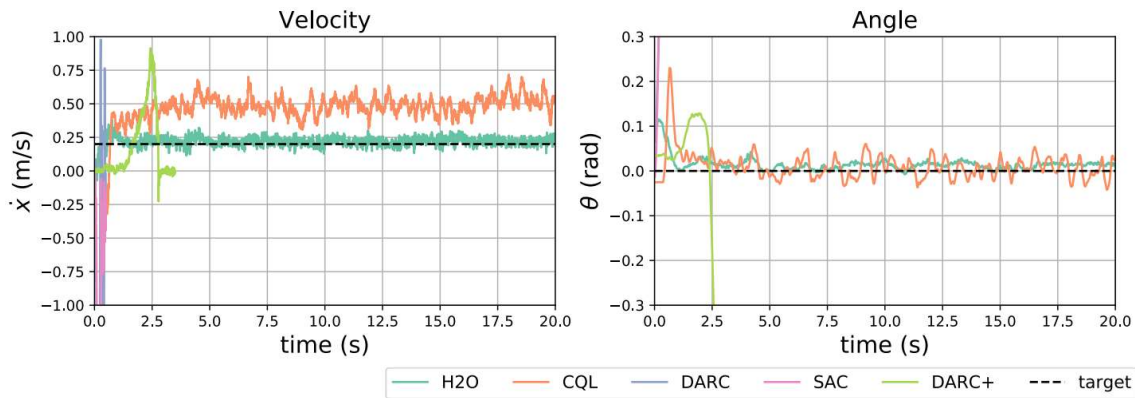
H2O



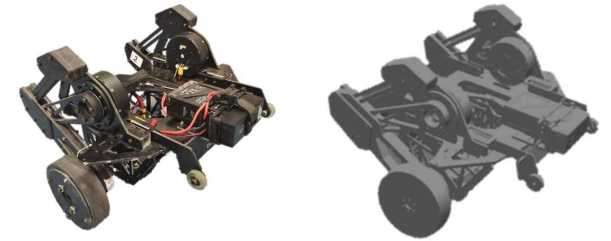


# Real-World Validation

## Moving Straight



## Wheel-Legged Robot



SAC

DARC

DARC+

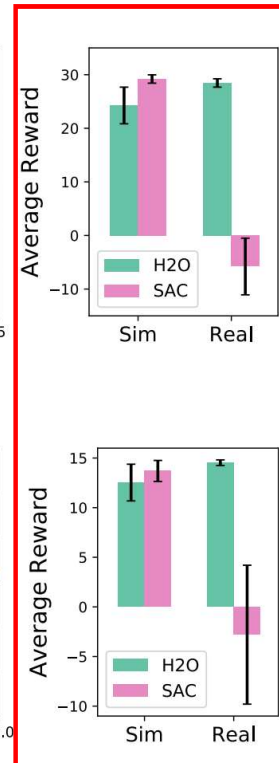
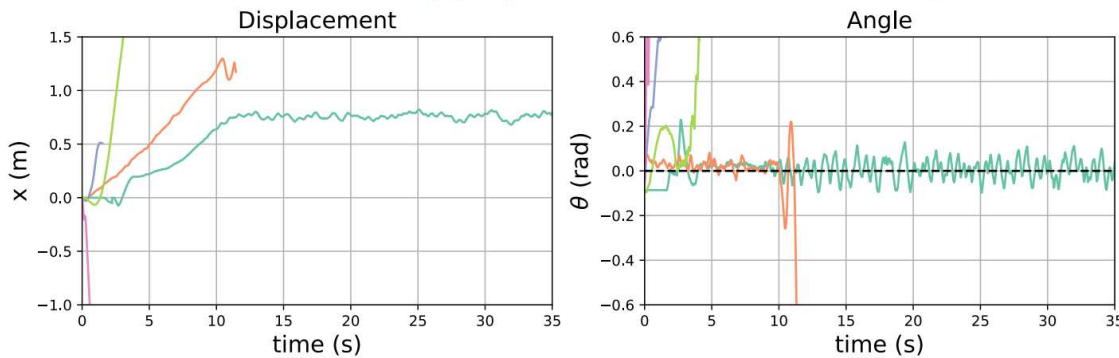
CQL

H2O

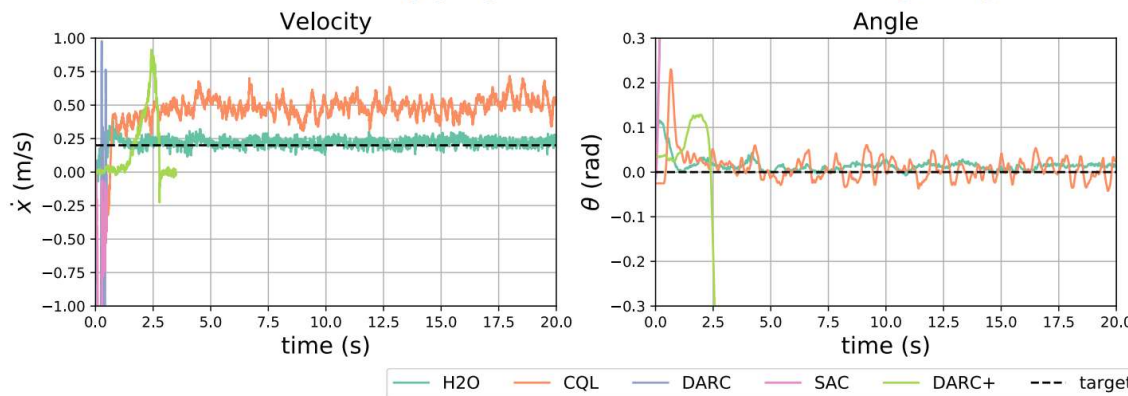


# Real-World Validation

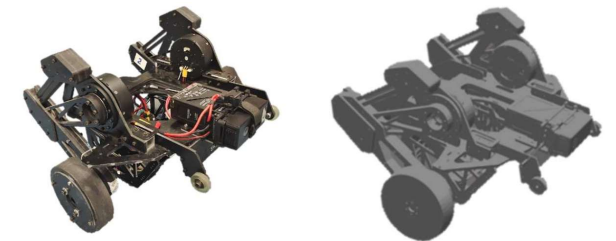
(c) Experimental results for standing still task



(d) Experimental results for moving straight task



## Wheel-legged Robot



**Rethinking: simulation-based evaluation can be very misleading, due to the dynamics gap.**

# When to Trust Your Simulator: Dynamics-Aware Hybrid Offline-and-Online Reinforcement Learning

Haoyi Niu, Shubham Sharma, Yiwen Qiu, Ming Li, Guyue Zhou, Jianming Hu, Xianyuan Zhan



## Thanks!



<https://github.com/t6-thu/H2O>



<https://www.youtube.com/watch?v=WRyEB6WEGc4>



[t6.da.thu@gmail.com](mailto:t6.da.thu@gmail.com)