# Diffusion-based Reinforcement Learning via Q-weighted Variational Policy Optimization

Shutong Ding[1,3]    Ke Hu[1]    Zhenhao Zhang[1]    Kan Ren[1,3]    Weinan Zhang[2]

Jingyi Yu[1,3]    Jingya Wang[1,3]    Ye Shi[1,3]

[1]ShanghaiTech University    [2]Shanghai Jiao Tong University
[3]MoE Key Laboratory of Intelligent Perception and Human Machine Collaboration
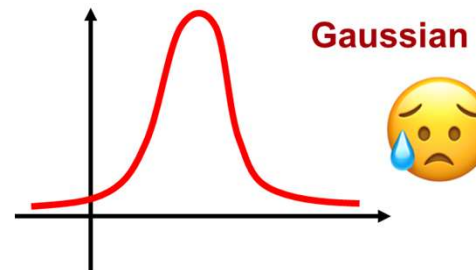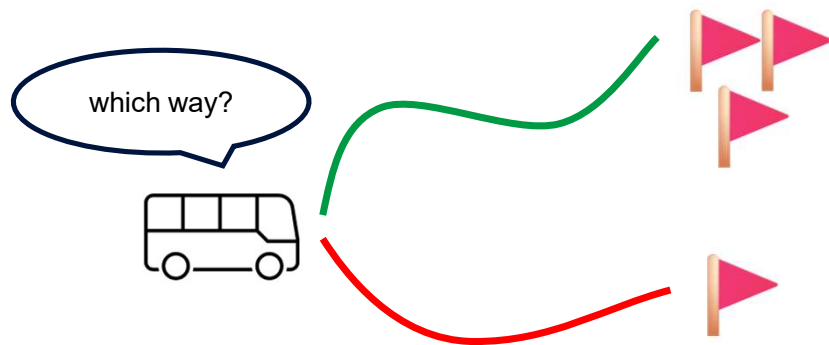
NeurIPS 2024

October 7, 2024
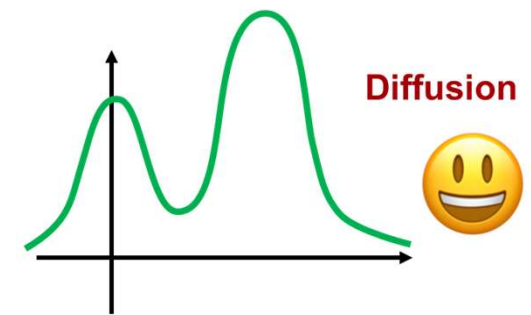
# Background: Diffusion in Online RL

1. **Exploration capability** of Gaussian policy or deterministic policy is limited

2. **Expressiveness and multimodality** of diffusion avoid policy falling into the local optimality

**Directions of applying diffusion in Online RL:**

1. **Use the variational loss of diffusion to do policy improvement**

$$\mathbb{E}_{t\sim[1,T],\mathbf{x}_0,\boldsymbol{\epsilon}_t} \left[||\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\boldsymbol{\epsilon}_t, t)||^2\right].$$

2. **Apply deterministic policy loss to train the diffusion model (like Diffusion-QL [1])**

$$\pi = \arg\min_{\pi_\theta}\mathcal{L}(\theta) = \mathcal{L}_d(\theta) + \mathcal{L}_q(\theta) = \mathcal{L}_d(\theta) - \boxed{\alpha \cdot \mathbb{E}_{\boldsymbol{s}\sim\mathcal{D},\boldsymbol{a}^0\sim\pi_\theta}\left[Q_\phi\left(\boldsymbol{s},\boldsymbol{a}^0\right)\right]}$$

[1] Wang Z, Hunt J J, Zhou M. Diffusion policies as an expressive policy class for offline reinforcement learning[J]. arXiv preprint arXiv:2208.06193, 2022..
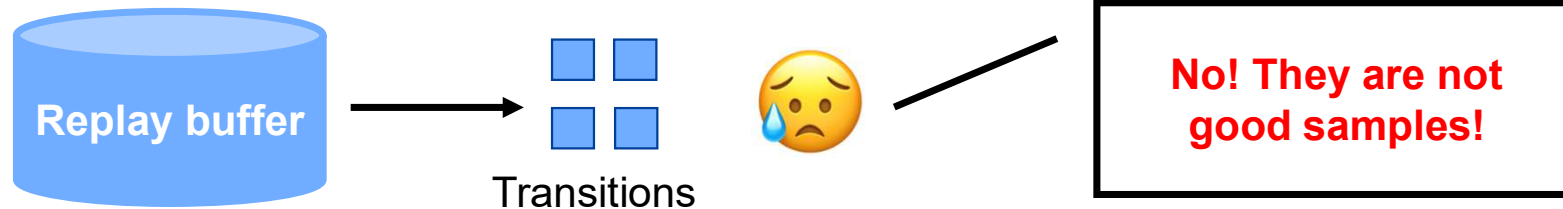
**Difficulties:**

**Direction 1:** Unavailability of **good** actions (not like imitation learning or offline RL)

Replay buffer

Transitions

No! They are not good samples!

**Direction 2: Too long** backpropagation chain (leads to high training cost and training instability)

$$a_T \quad \dots \quad a_t \quad a_{t-1} \quad \dots \quad a_0$$

**Backpropagation chain**

Too long chain…

## Previous works:

**DIPO [1]:** Directly perform gradient update on action sample (**affect the multimodality**)

$$\mathbf{a}_t + \eta \nabla_{\mathbf{a}} Q_\pi(\mathbf{s}_t, \mathbf{a}_t) \rightarrow \mathbf{a}_t$$

$$\pi \xrightarrow{\text{data}} \mathcal{D} = \{\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}, r_{t+1}\} \xrightarrow{\text{action gradient}} \mathcal{D}' = \{\mathbf{s}_t, \mathbf{a}_t\} \xrightarrow{\text{diffusion policy}} \pi'$$
$$\pi \leftarrow \pi'$$

**QSM [2]:** Do score matching with the gradient of Q function (**Doubled approximation error**)

Update score model:
$$\phi = \operatorname{argmin}_\phi N^{-1} \sum (\Psi_\phi(x_t, a_t) - \nabla_a Q(x_t, a_t))^2;$$

[1] Yang L, Huang Z, Lei F, et al. Policy representation via diffusion probability model for reinforcement learning[J]. arXiv preprint arXiv:2305.13122, 2023.
[2] Psenka M, Escontrela A, Abbeel P, et al. Learning a diffusion model policy from rewards via q-score matching[J]. arXiv preprint arXiv:2312.11752, 2023.

# Solution: Q-weighted Variational Policy Optimization

上海科技大学
ShanghaiTech University

Can we use the variational loss to train diffusion policy?

Distinguish bad and good samples according to their **Q-value**?

But samples from replay buffer are still not good enough!
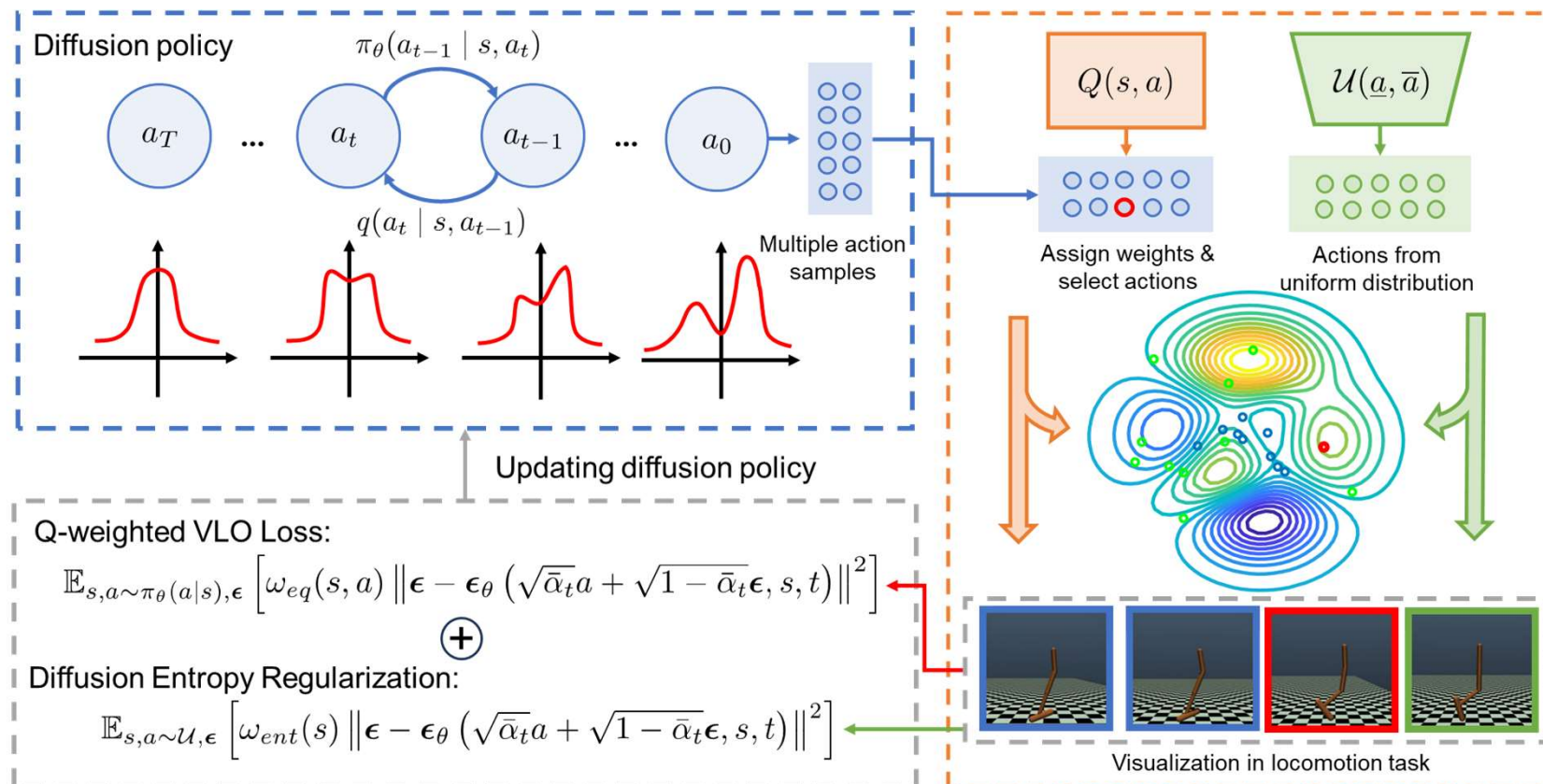
What if use samples from **current policy**?

# Solution: Q-weighted Variational Policy Optimization

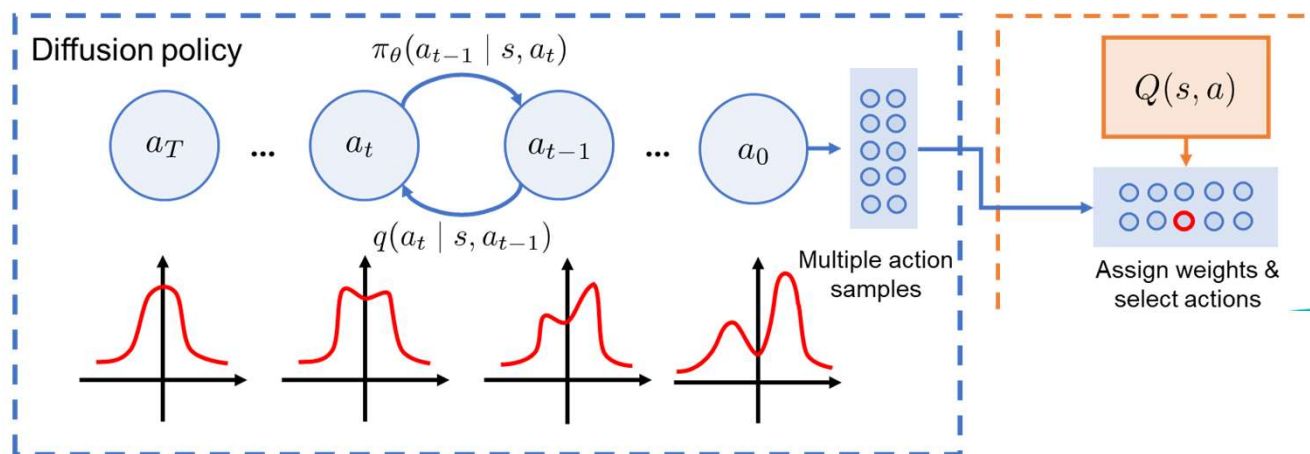Motivated by these two ideas, we propose Q-weighted variational policy optimization. 😃

Diffusion policy

$\pi_\theta(a_{t-1} \mid s, a_t)$

$a_T$ ... $a_t$ $a_{t-1}$ ... $a_0$

$q(a_t \mid s, a_{t-1})$

Multiple action samples

Updating diffusion policy

Q-weighted VLO Loss:

$$\mathbb{E}_{s,a\sim\pi_\theta(a\mid s),\epsilon}\left[\omega_{eq}(s,a)\left\|\epsilon - \epsilon_\theta\left(\sqrt{\bar{\alpha}_t}a + \sqrt{1-\bar{\alpha}_t}\epsilon, s, t\right)\right\|^2\right]$$

$\oplus$

Diffusion Entropy Regularization:

$$\mathbb{E}_{s,a\sim\mathcal{U},\epsilon}\left[\omega_{ent}(s)\left\|\epsilon - \epsilon_\theta\left(\sqrt{\bar{\alpha}_t}a + \sqrt{1-\bar{\alpha}_t}\epsilon, s, t\right)\right\|^2\right]$$

$Q(s,a)$  $\mathcal{U}(\underline{a},\bar{a})$

Assign weights & select actions

Actions from uniform distribution

Visualization in locomotion task

上海科技大学
ShanghaiTech University
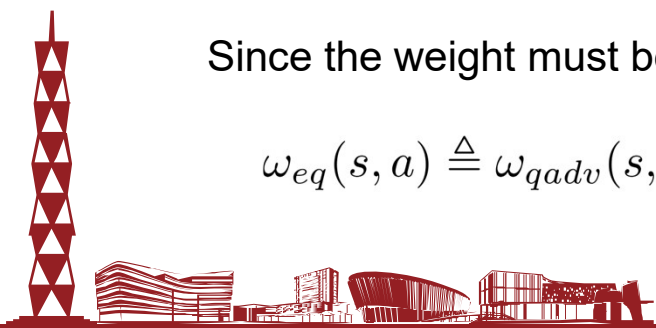
# Q-weighted Variational Loss



Q-weighted VLO Loss (tight lower bound of RL policy objective):

$$\mathbb{E}_{s,a\sim\pi_\theta(a|s),\boldsymbol{\epsilon}}\left[\omega_{eq}(s,a)\left\|\boldsymbol{\epsilon}-\boldsymbol{\epsilon}_\theta\left(\sqrt{\bar{\alpha}_t}a+\sqrt{1-\bar{\alpha}_t}\boldsymbol{\epsilon},s,t\right)\right\|^2\right]$$

Since the weight must be nonnegative, we define the weight as

$$\omega_{eq}(s,a)\triangleq\omega_{qadv}(s,a)=\begin{cases} A(s,a), & A(s,a)\geq 0 \\ 0, & A(s,a)<0 \end{cases},$$

**How to Further improve the training action quality?**

**Only choose the best generated sample for training.**

# Diffusion Entropy Regularization
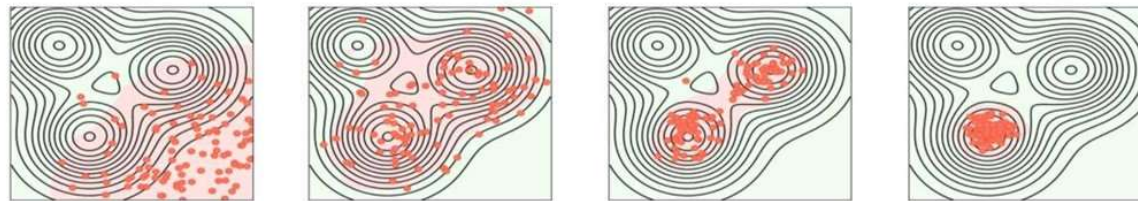
We also add Diffusion Entropy Regularization term in objective for enhancing exploration capability:

$$\mathbb{E}_{s,a\sim\mathcal{U},\boldsymbol{\epsilon}}\left[\omega_{ent}(s)\left\|\boldsymbol{\epsilon}-\boldsymbol{\epsilon}_\theta\left(\sqrt{\bar{\alpha}_t}a+\sqrt{1-\bar{\alpha}_t}\boldsymbol{\epsilon},s,t\right)\right\|^2\right]$$

- Maximizing entropy can be viewed as **minimizing the distance** between current policy and the **uniform distribution**.

- Hence, we use samples from the uniform distribution to train diffusion policy for maximizing the diffusion entropy.



$Q(s,a)$

$\mathcal{U}(\underline{a},\overline{a})$

Assign weights & select actions

Actions from uniform distribution

We also add Diffusion Entropy Regularization term in objective for enhancing exploration capability:

$$\mathbb{E}_{s,a\sim\mathcal{U},\boldsymbol{\epsilon}}\left[\omega_{ent}(s)\left\|\boldsymbol{\epsilon}-\boldsymbol{\epsilon}_\theta\left(\sqrt{\bar{\alpha}_t}a+\sqrt{1-\bar{\alpha}_t}\boldsymbol{\epsilon},s,t\right)\right\|^2\right]$$
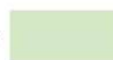


Diffusion Policy w/o entropy term

Training Procedure

Diffusion Policy w/ entropy term
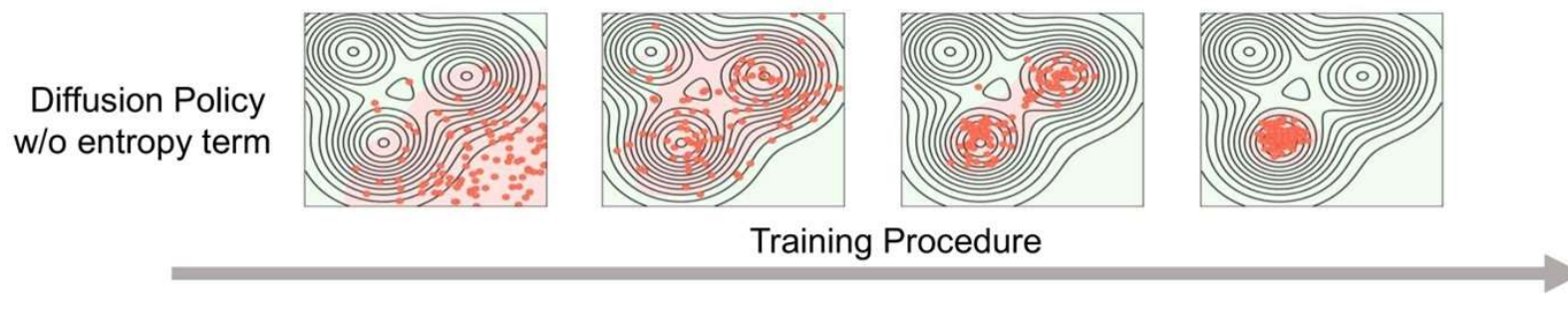
Explorable area of diffusion policy    Unexplorable area    Target objective function

上海科技大学
ShanghaiTech University

- Diffusion policy has a **large policy variance**

- Reducing this variance can improve **the quality** of collected transitions

Diffusion Policy
w/o entropy term

Training Procedure

We design **efficient behavior policy** via action selection for sample efficiency:

$$\pi_\theta^K(a \mid s) \triangleq \underset{a \in \{a_1, \cdots, a_K \sim \pi_\theta(a|s)\}}{\mathrm{argmax}} Q(s, a).$$
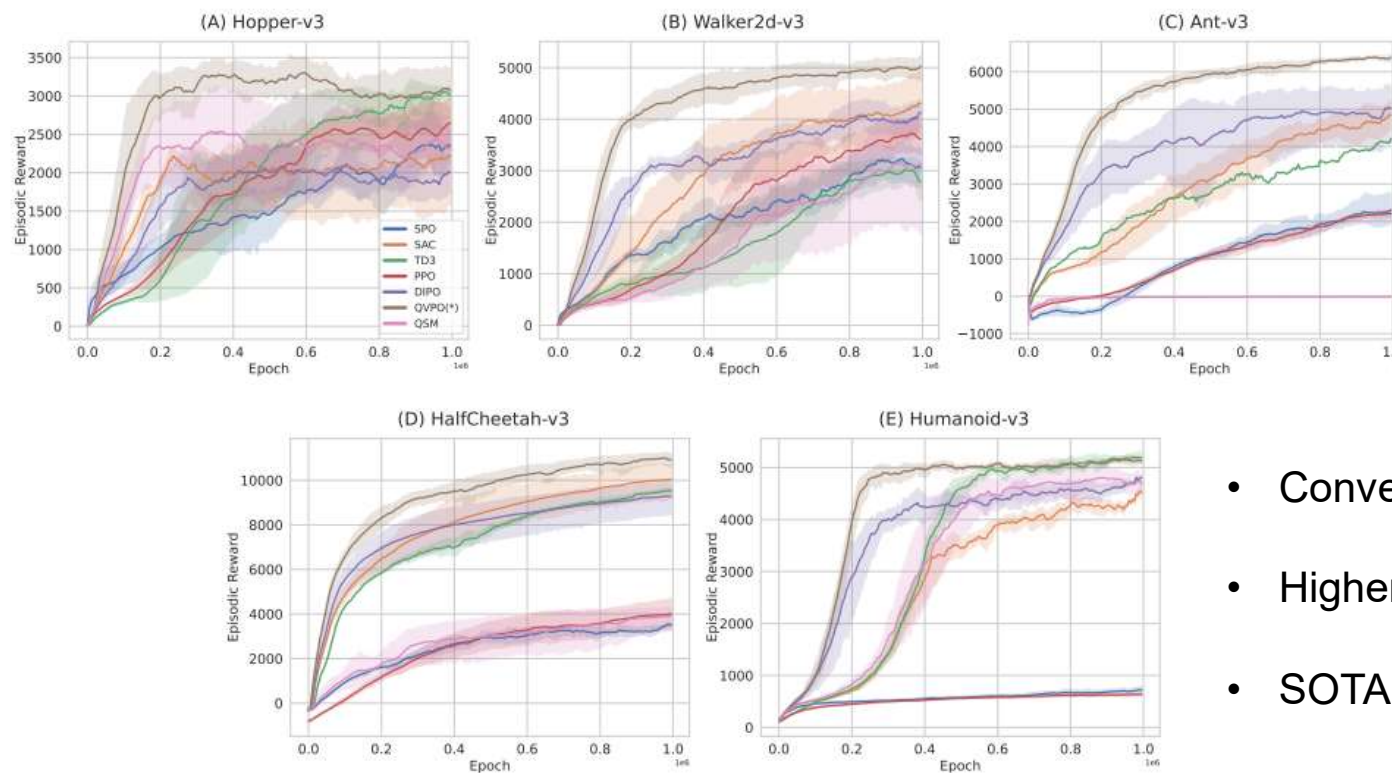
立志成才报国裕民

# Results



(A) Hopper-v3

(B) Walker2d-v3

(C) Ant-v3

(D) HalfCheetah-v3

(E) Humanoid-v3

- Converge faster (sample efficiency)

- Higher cumulative reward

- SOTA method in online RL

# Comparison with SAC in Humanoid-v3

SAC

our QVPO
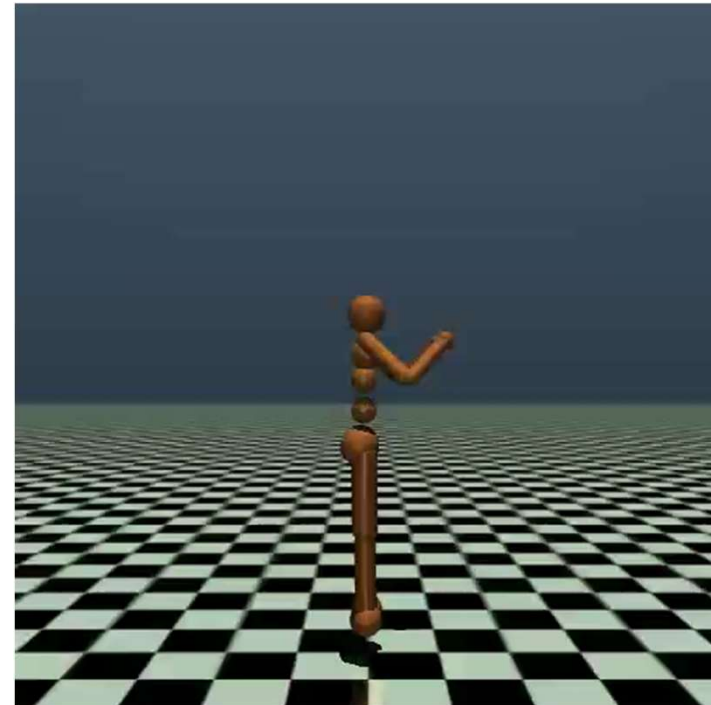
# More Tasks in Unitree Humanoid H1

Balance (TD-MPC2)

Balance (our QVPO)

# More Tasks in Unitree Humanoid H1

Walk (QVPO)
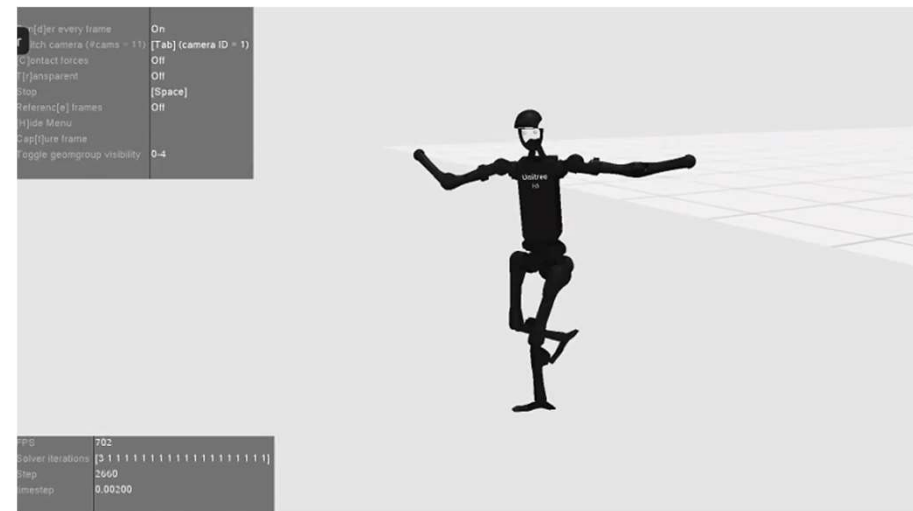
Run (QVPO)