



# Maximum Entropy Reinforcement Learning via Energy-Based Normalizing Flow

Chen-Hao Chao<sup>\*1 2</sup>, Chien Feng<sup>\*1</sup>, Wei-Fang Sun<sup>2</sup>, Cheng-Kuang Lee<sup>2</sup>, Simon See<sup>2</sup>, and Chun-Yi Lee<sup>1</sup>  
(\* equal contribution)

<sup>1</sup> Elsa Lab, National Tsing Hua University, Hsinchu City, Taiwan  
<sup>2</sup> NVIDIA AI Technology Center, NVIDIA Corporation, Santa Clara, CA, USA



## TL'DR

We introduce a novel Maximum Entropy (MaxEnt) Reinforcement Learning (RL) framework that supports:

- **Single training loss optimization process**
- **Exact soft value calculation**
- **Multi-modal action distribution modeling**

Our method achieves superior performance on:

- **the MuJoCo benchmark**
- **the Omniverse Isaac Gym Environments**

compared to widely-adopted baselines (e.g., SAC [1]).



## Background

### Maximum Entropy Reinforcement Learning

MaxEnt RL augments the standard RL objective with the entropy of a policy as follows:

$$\pi_{\text{MaxEnt}}^* = \arg\max_{\pi} \sum_t \mathbb{E}_{\rho_{\pi}}[r_t + \alpha \mathcal{H}(\pi(\cdot | s_t))], \quad (1)$$

where  $\mathcal{H}$  is the entropy and  $\alpha$  is a temperature parameter. The solution  $\pi_{\text{MaxEnt}}^*$  can be expressed as:

$$\pi_{\text{MaxEnt}}^*(a_t | s_t) = \exp\left(\frac{1}{\alpha} (Q_{\text{soft}}^*(s_t, a_t) - V_{\text{soft}}^*(s_t))\right), \quad (2)$$

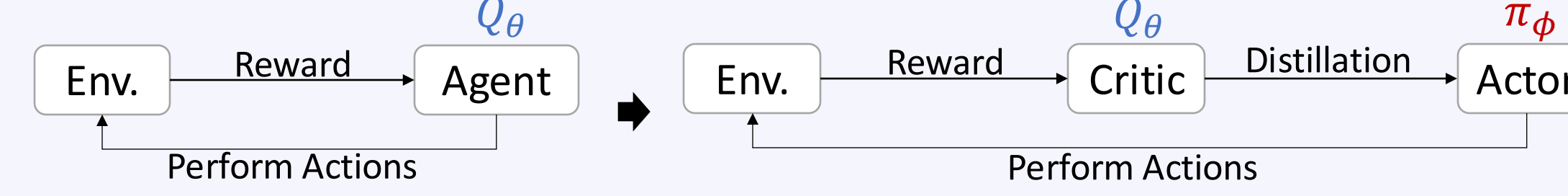
where  $Q_{\text{soft}}^*(s_t, a_t)$  is the optimal soft Q-function and  $V_{\text{soft}}^*(s_t) = \alpha \log \int \exp(\frac{1}{\alpha} Q_{\text{soft}}^*(s_t, a)) da$ . Given an experience reply buffer  $\mathcal{D}$ , the objective of  $Q_{\theta}$  is written as:

$$\mathcal{L}(\theta) = \mathbb{E}_{\mathcal{D}} \left[ \frac{1}{2} \left( Q_{\theta}(s_t, a_t) - (r_t + \gamma V_{\theta}(s_{t+1})) \right)^2 \right]. \quad (3)$$

### Actor-Critic Designs

- **Two Key Challenges of MaxEnt RL with Conti.  $a_t$** 
  - **Inefficient sampling** using  $\pi_{\theta}(a_t | s_t) \propto \exp(Q_{\theta}(s_t, a_t))$
  - **Intractable integration** operation in  $V_{\theta}(s_t)$

- **Actor-Critic Frameworks [1-5]**



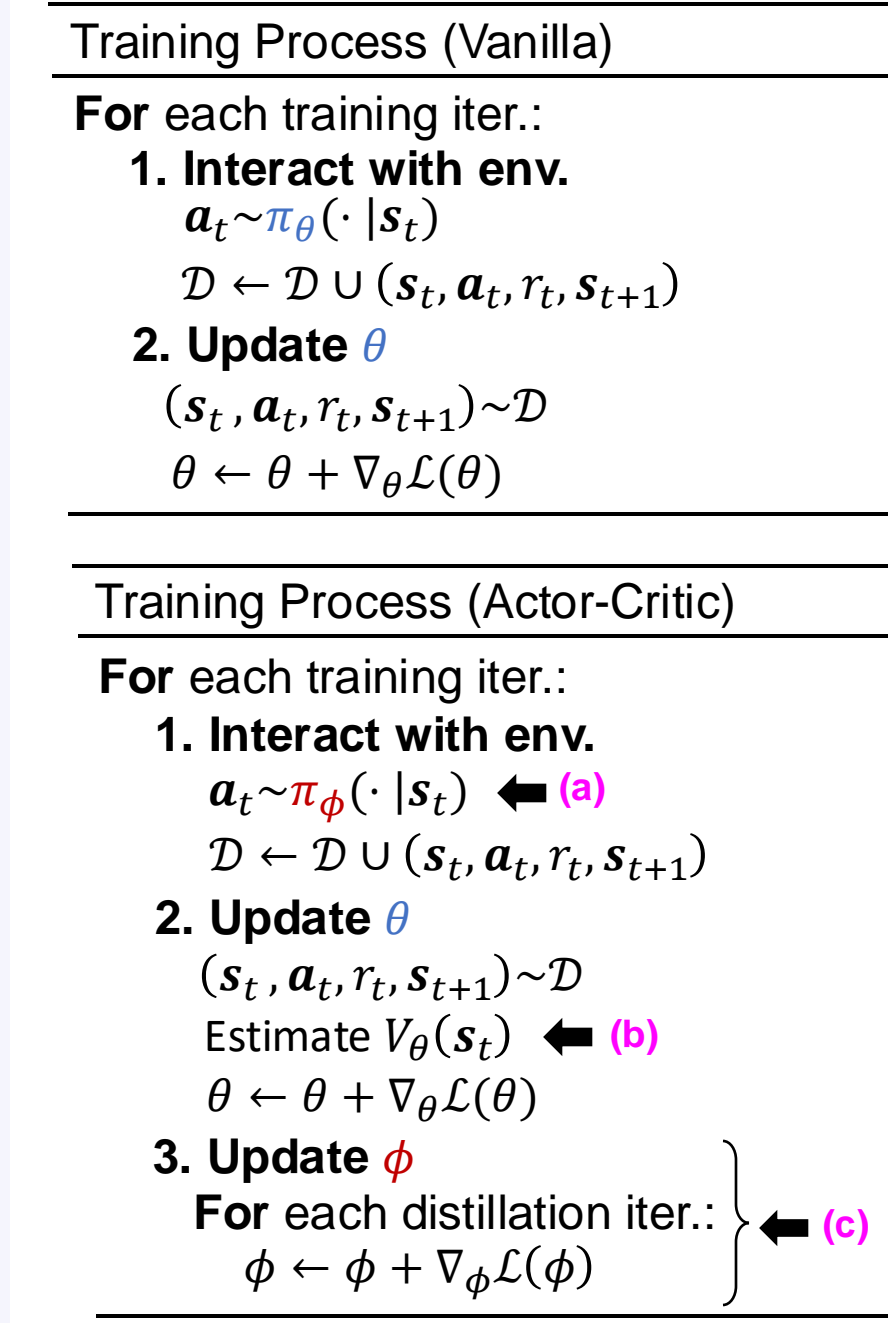
- Actor Loss (Reverse KL Divergence  $\mathbb{D}_{KL}[\pi_{\phi} || \pi_{\theta}]$ ):

$$\mathcal{L}(\phi) = \mathbb{E}_{\mathcal{D}, \pi_{\phi}}[-(Q_{\theta}(s_t, a_t) - \alpha \log \pi_{\phi}(a_t | s_t))] \quad (4)$$

- Soft Value Estimation (Monte Carlo Estimation):

$$\text{Soft Q-Learning (SQL) [1]: } V_{\theta, \phi}(s_t) \approx \alpha \log \frac{1}{M} \sum_{i=1}^M \left( \frac{\exp(Q_{\theta}(s_t, a^i)/\alpha)}{\pi_{\phi}(a^i | s_t)} \right) \quad (5)$$

$$\text{Soft Actor-Critic (SAC) [2]: } V_{\theta, \phi}(s_t) \approx \frac{1}{M} \sum_{i=1}^M (Q_{\theta}(s_t, a^i) - \alpha \log \pi_{\phi}(a^i | s_t)) \quad (6)$$



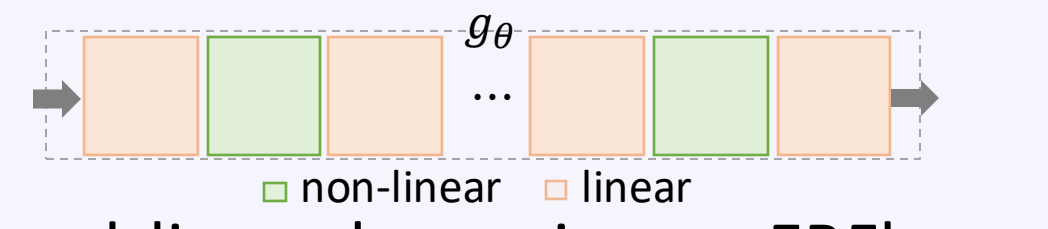
### Energy-based Normalizing Flows

- **Normalizing Flows**

The probability density functions (pdf)  $p_{\theta}$  are parameterized using a prior distribution  $p_{\mathbf{u}}$  of a variable  $\mathbf{u}$  and an invertible mapping  $g = g_L \circ \dots \circ g_1$ , where  $g_{\theta}^i: \mathbb{R}^D \rightarrow \mathbb{R}^D, i \in \{1, \dots, L\}$ . Based on the change of variable theorem,  $p_{\theta}$  can be expressed as:

$$p_{\theta}(\mathbf{x}) = p_{\mathbf{u}}(g_{\theta}(\mathbf{x})) \prod_{i=1}^L |\det(\mathcal{J}_{g_{\theta}^i}(\mathbf{x}^{i-1}))|, \quad (7)$$

where  $\mathcal{J}_{g_i}$  represents the Jacobian of  $g_i$ .



- **Energy-based Normalizing Flows (EBFlow) [6]**

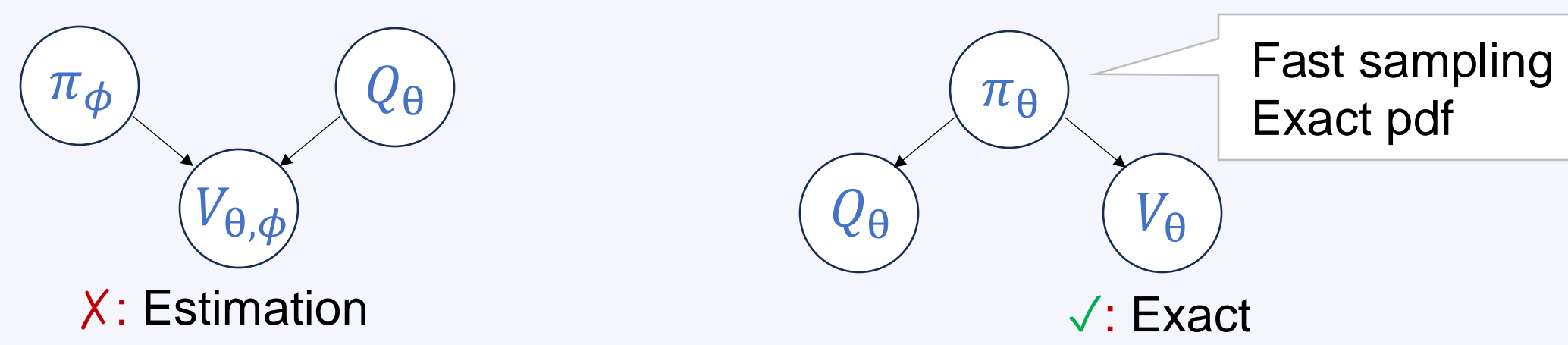
Let  $\mathcal{S}_n$  and  $\mathcal{S}_l$  be the index sets of non-linear and linear layers in  $g_{\theta}$ . EBFlow reinterprets Eq. (7) as a Boltzmann distribution and factorizes it as follows:

$$p_{\theta}(\mathbf{x}) = p_{\mathbf{u}}(g_{\theta}(\mathbf{x})) \underbrace{\prod_{i \in \mathcal{S}_n} |\det(\mathcal{J}_{g_{\theta}^i}(\mathbf{x}^{i-1}))|}_{\triangleq \exp(-E_{\theta}(\mathbf{x}))} \underbrace{\prod_{i \in \mathcal{S}_l} |\det(\mathcal{J}_{g_{\theta}^i})|}_{\triangleq Z_{\theta}^{-1}} \quad (8)$$

## Methodology

### MaxEnt RL via EBFlow (MEow) Framework

- **Previous:** Model  $Q_{\theta}$  and  $\pi_{\phi}$  and derive  $V_{\theta, \phi}$
- **Ours:** Model  $\pi_{\theta}$  and extract  $Q_{\theta}$  and  $V_{\theta}$



The MEow framework parameterizes the policy as a state-conditioned EBFlow:

$$\pi_{\theta}(a_t | s_t) = p_{\mathbf{u}}(g_{\theta}(a_t | s_t)) \prod_{i \in \mathcal{S}_n} |\det(\mathcal{J}_{g_{\theta}^i}(a_t^{i-1} | s_t))| \prod_{i \in \mathcal{S}_l} |\det(\mathcal{J}_{g_{\theta}^i}(s_t))| \quad (9)$$

$$\triangleq \exp\left(\frac{1}{\alpha} Q_{\theta}(s_t, a_t)\right) \triangleq \exp\left(-\frac{1}{\alpha} V_{\theta}(s_t)\right)$$

**Proposition 1.** Eq. (9) satisfies the following statements:

- Given that the Jacobian of  $g_{\theta}$  is non-singular,  $V_{\theta}(s) \in \mathbb{R}$  and  $Q_{\theta}(s, a) \in \mathbb{R}$ .
- $V_{\theta}(s) = \alpha \log \int \exp(\frac{1}{\alpha} Q_{\theta}(s, a)) da$ .

### Expressivity

- (✓) Normalizing flows are universal approximators [7]
- (✓)  $Q_{\theta}$  and  $V_{\theta}$  can represent any real number (i.e., Proposition 1)

### Training

- (X  $\rightarrow$  ✓) Single training loss optimization process (i.e., Eq. (3))
- (X  $\rightarrow$  ✓) Exact soft value calculation (i.e., Eq. (9))

### Inference

- (X  $\rightarrow$  ✓) Consistent policy for training and inference (i.e.,  $\pi_{\theta}$ )

### Deterministic Inference Technique of MEow

**Proposition 2.** Given that the Jacobians are constants w.r.t. the input, then:

$$\arg\max_{a_t} Q_{\theta}(s_t, a_t) = g_{\theta}^{-1}(\arg\max_{\mathbf{u}} p_{\mathbf{u}}(\mathbf{u}) | s_t). \quad (10)$$

### Practical Implementation

- Adopt NICE-like [10] architecture for  $\pi_{\theta}$
- Using the mean  $\mu$  of Gaussian (i.e.,  $g_{\theta}^{-1}(\mu | s_t)$ ) for inference.

## Experiments

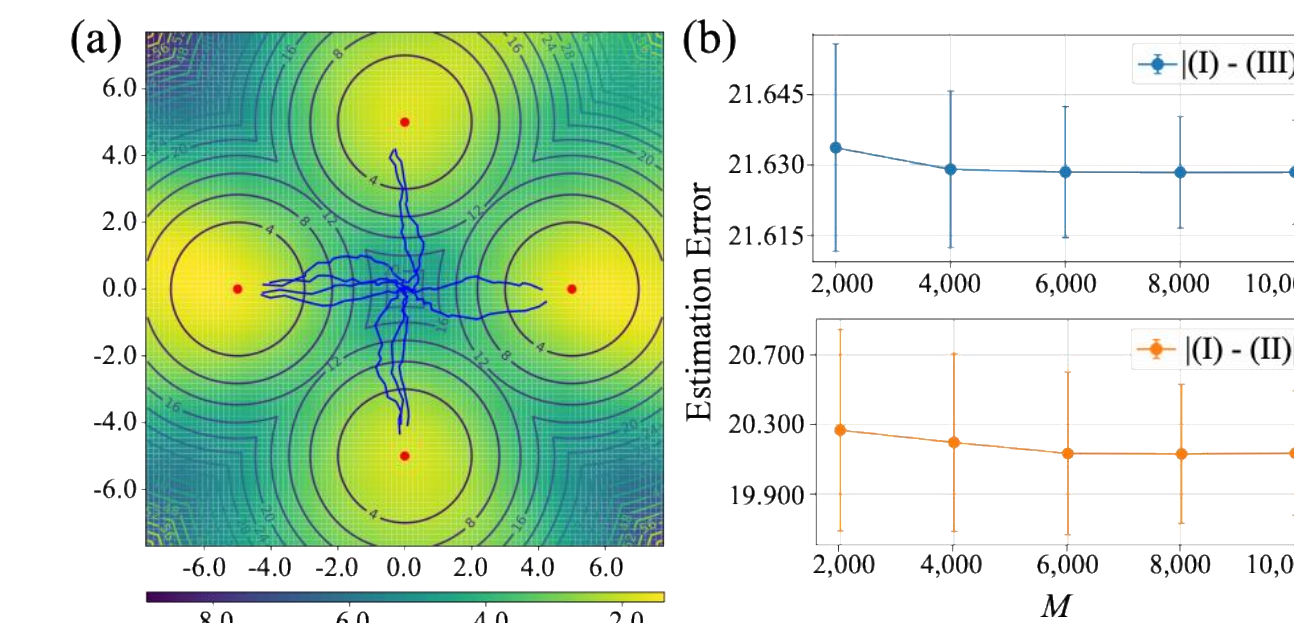
### Setups

- **Environments:** Multi-goals [1], MuJoCo [8], Isaac [9]
- **Architecture:**
  - Additive coupling layers [10]
  - Element-wise linear layer

### Baselines

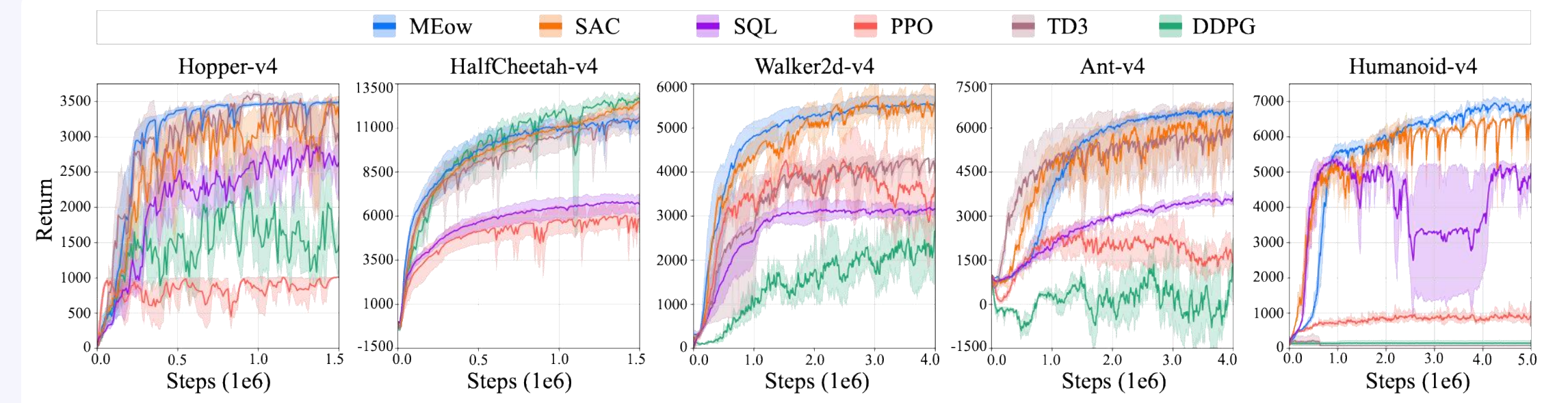
- SAC
- SQL
- PPO
- DDPG
- TD3

### Multi-Goal Environment



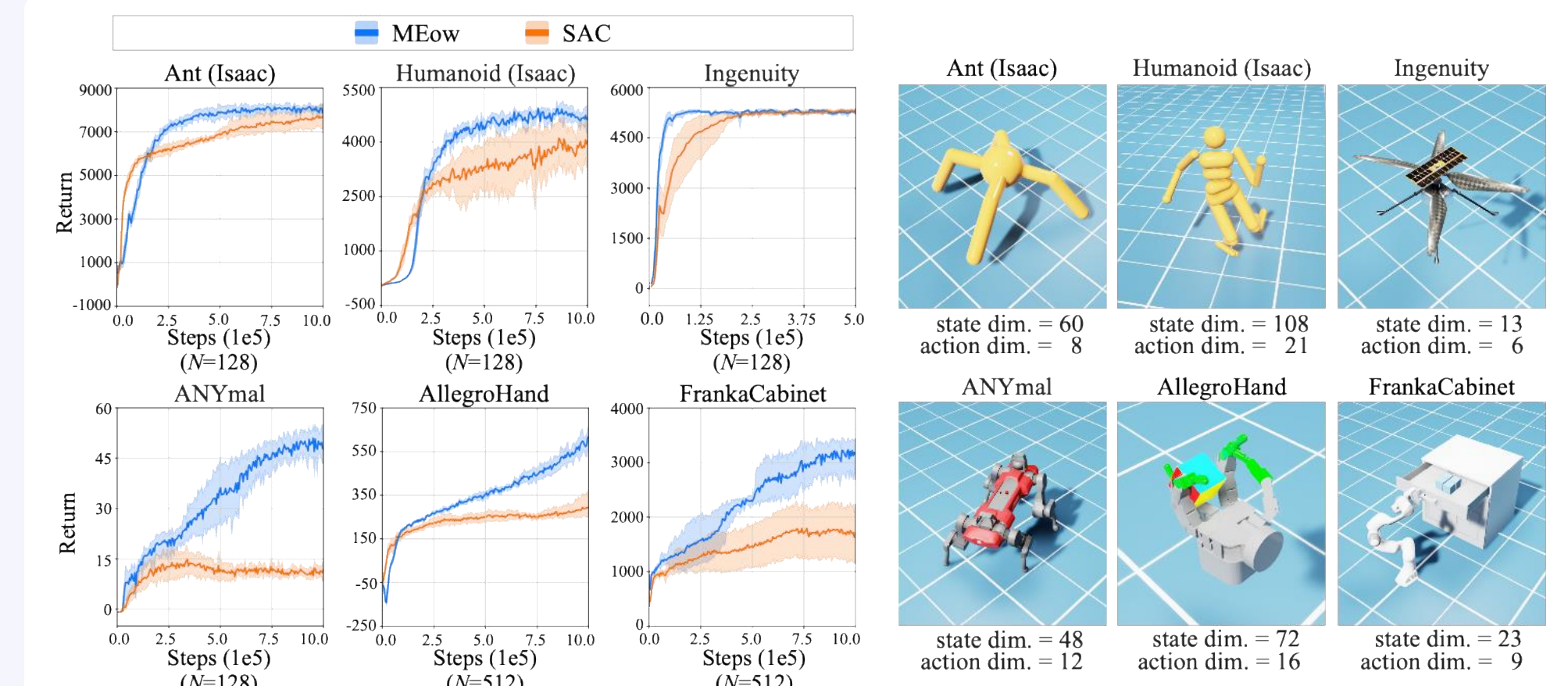
**Figure 1.** (a) The soft value function and the trajectories generated using our method on the multi-goal environment. (b) The estimation error evaluated at the initial state under different choices of  $M$ . In this plot, (I), (II), and (III) represent  $V_{\theta}$  in Eq. (9), Eq. (5), and Eq. (6), respectively.

### MuJoCo Benchmark



**Figure 2.** The results in terms of total returns versus the number of training steps evaluated on five MuJoCo environments. Each curve represents the mean performance, with shaded areas indicating the 95% confidence intervals, derived from five independent runs with different seeds.

### Omniverse Isaac Gym Environments



**Figure 3 (Left).** A comparison of SAC and MEow on six Isaac Gym environments. Each curve represents the mean performance of five runs, with shaded areas indicating the 95% confidence intervals. 'Steps' in the x-axis represents the number of training steps, each of which consists of  $N$  parallelizable interactions with the environments.

**Figure 3 (Right).** A demonstration of six Isaac Gym environments. The dimensionalities of the state and action for each environment are denoted below each subfigure.

## References

- [1] Haarnoja *et al.* Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor, *ICML* 2018.
- [2] Haarnoja *et al.* Reinforcement Learning with Deep Energy-Based Policies, *ICML* 2017.
- [3] Haarnoja *et al.* Latent Space Policies for Hierarchical Reinforcement Learning, *ICML* 2018.

- [4] Zhang *et al.* Latent State Marginalization as a Low-cost Approach for Improving Exploration, *ICLR* 2023.
- [5] Messaoud *et al.* S2AC: Energy-Based Reinforcement Learning with Stein Soft Actor Critic, *ICLR* 2024.
- [6] Chao *et al.* Training Energy-Based Normalizing Flow with Score-Matching Objectives, *NeurIPS* 2023.
- [7] Papamakarios *et al.* Normalizing Flows for Probabilistic Modeling and Inference. *JMLR* 2019.

- [8] Todorov *et al.* MuJoCo: A physics engine for model-based control. *IROS* 2012.
- [9] Makoviychuk *et al.* Isaac Gym: High-Performance GPU-Based Physics Simulation For Robot Learning. 2021.
- [10] Dinh *et al.* NICE: Non-linear Independent Components Estimation. *ICLR Workshop* 2015.

## Acknowledgements

