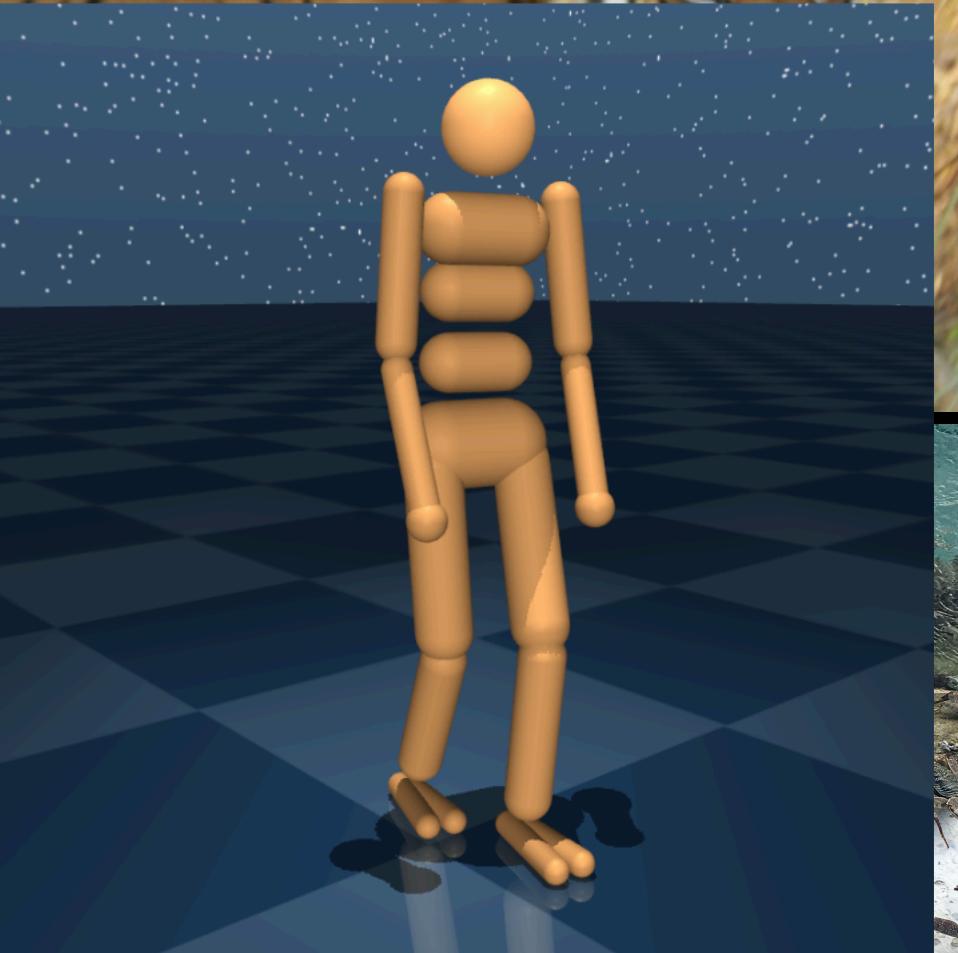


# **reinforcement learning for locomotion**

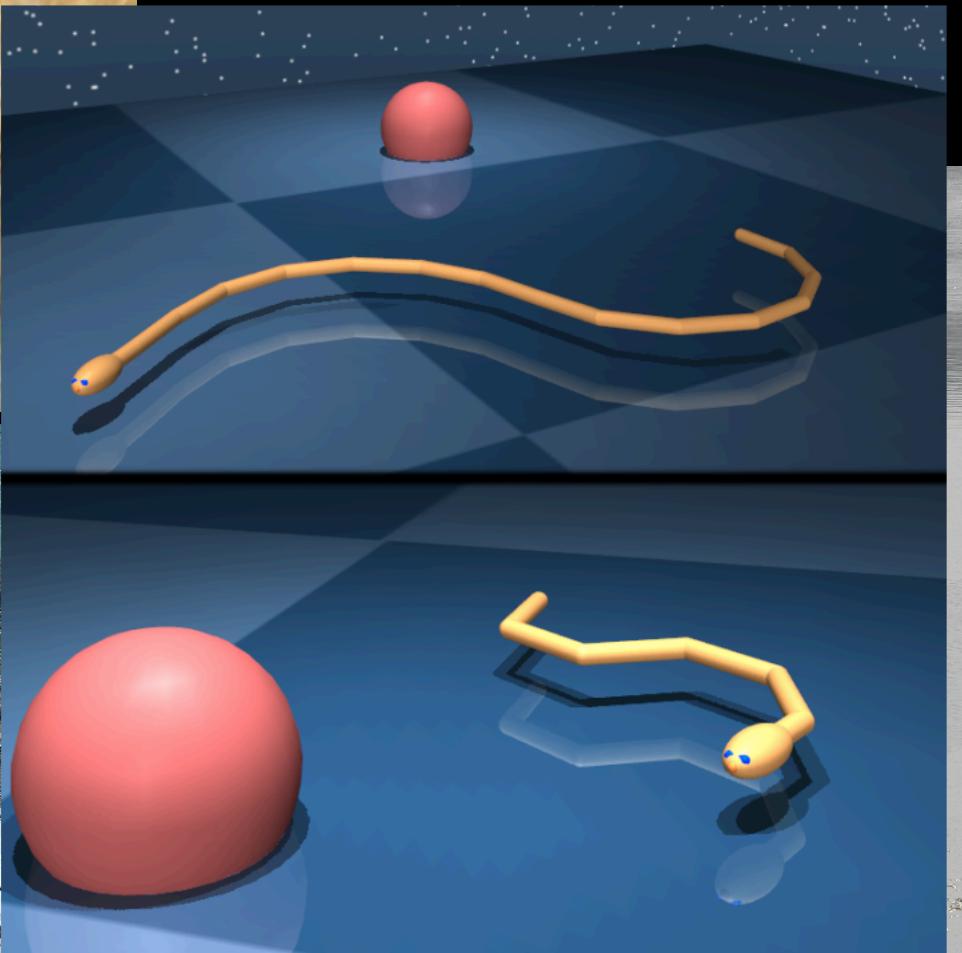
**Shruti Mishra, Abbas Abdolmaleki,  
Arthur Guez, Piotr Trochim, Doina Precup**

reinforcement learning  
biologically-inspired domain

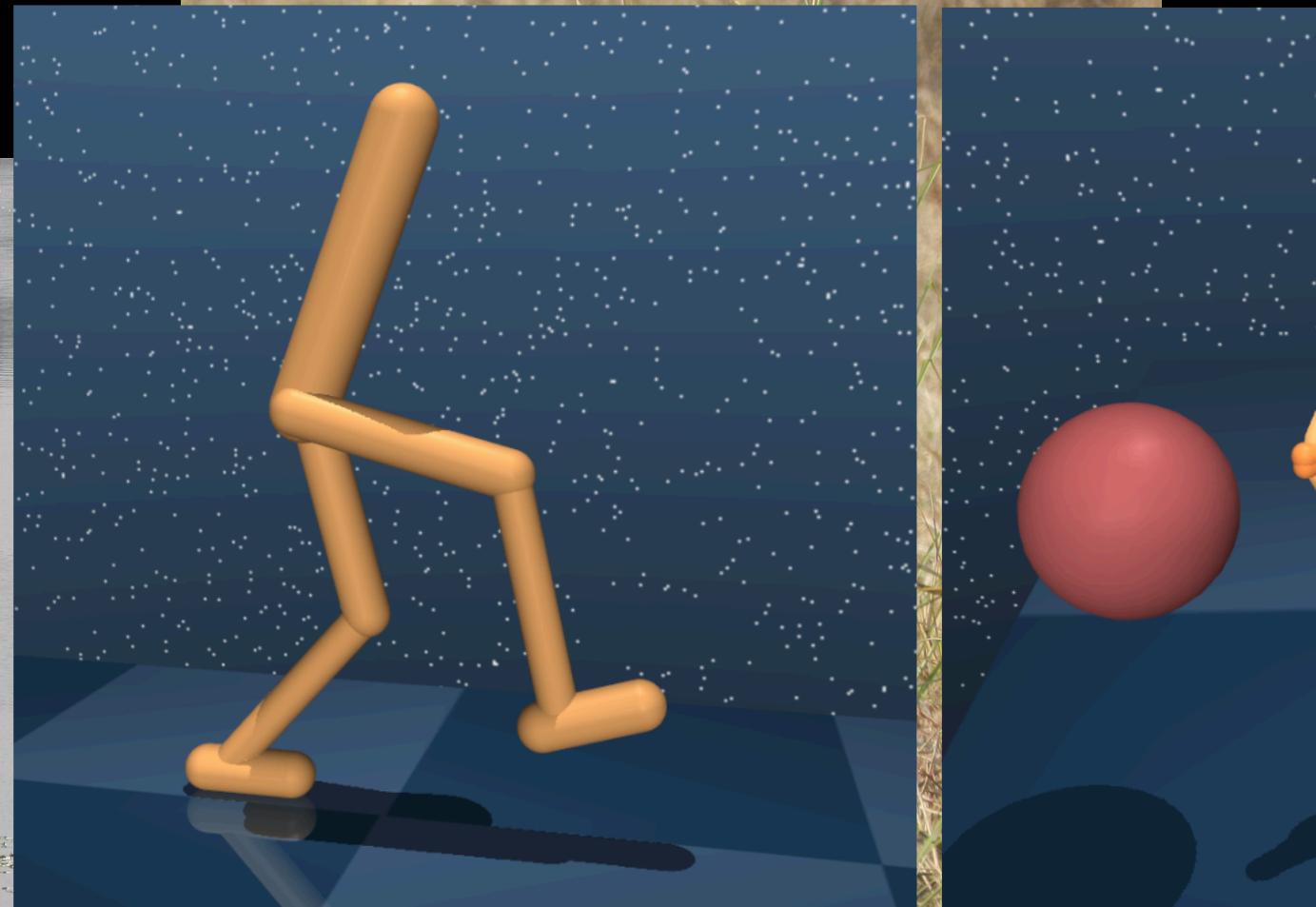
# symmetry biologically-inspired domain



Humanoid



Swimmer



Planar Walker



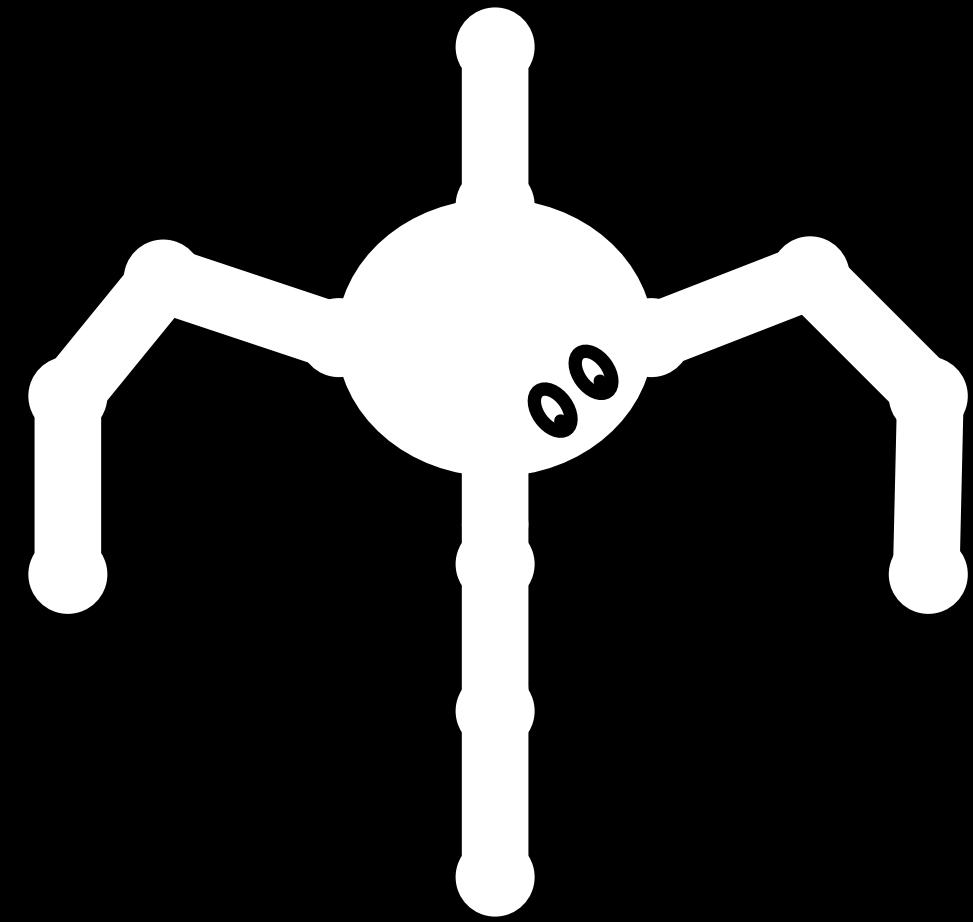
Fish



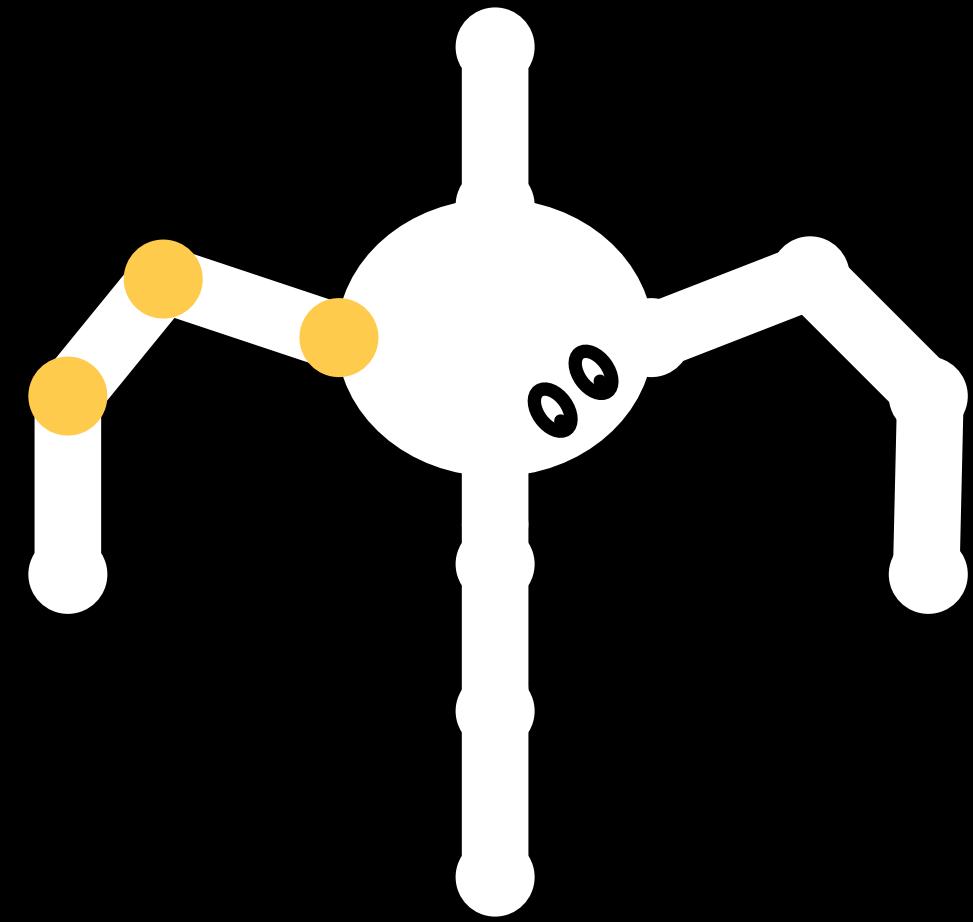
Augmenting learning using **symmetry**  
in a **biologically-inspired domain**

Augmenting learning using symmetry  
in a biologically-inspired domain

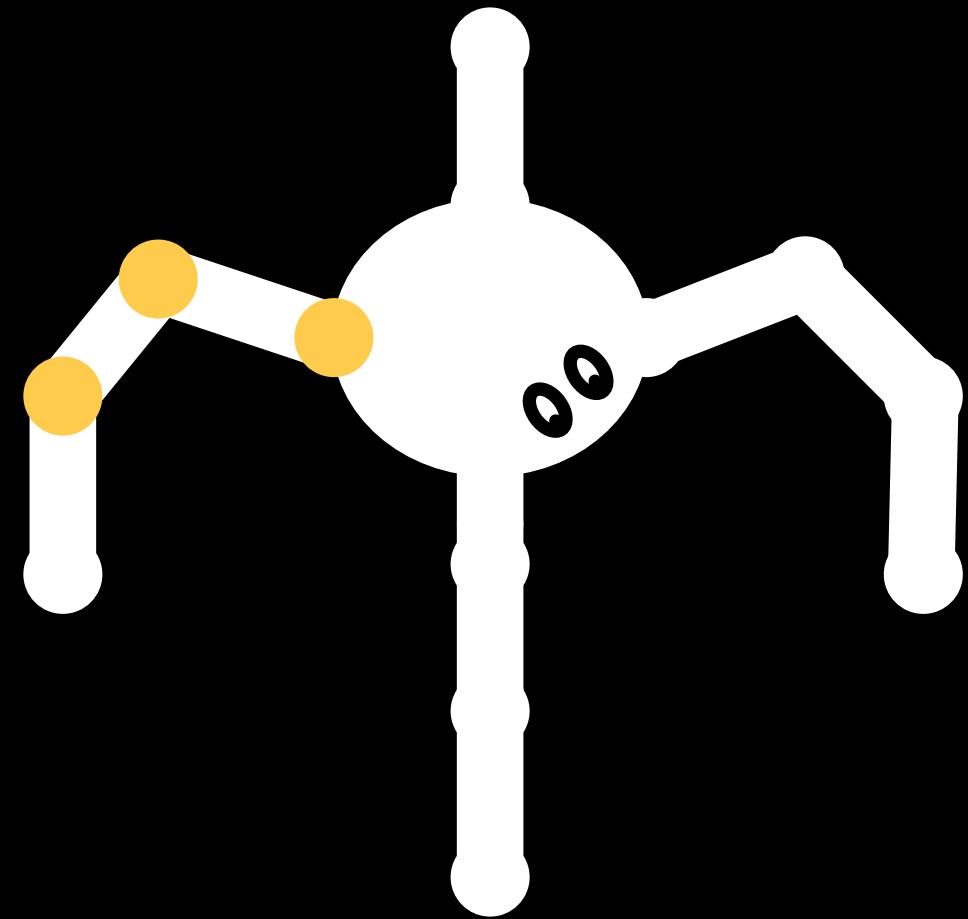
# Domain: Quadruped



# Domain: Quadruped



# Domain: Quadruped



12 Actions

Each leg has 3 hinge joints → 12

78 Observations

Torso: Orientation relative to horizontal → 1

Torso: Forward velocity → 3

Torso: Accelerometer and gyroscope → 6

Each toe: Force and Torque → 12

Each joint: Instantaneous angle, angular velocity, activation → 44

# Task: Walk and Run

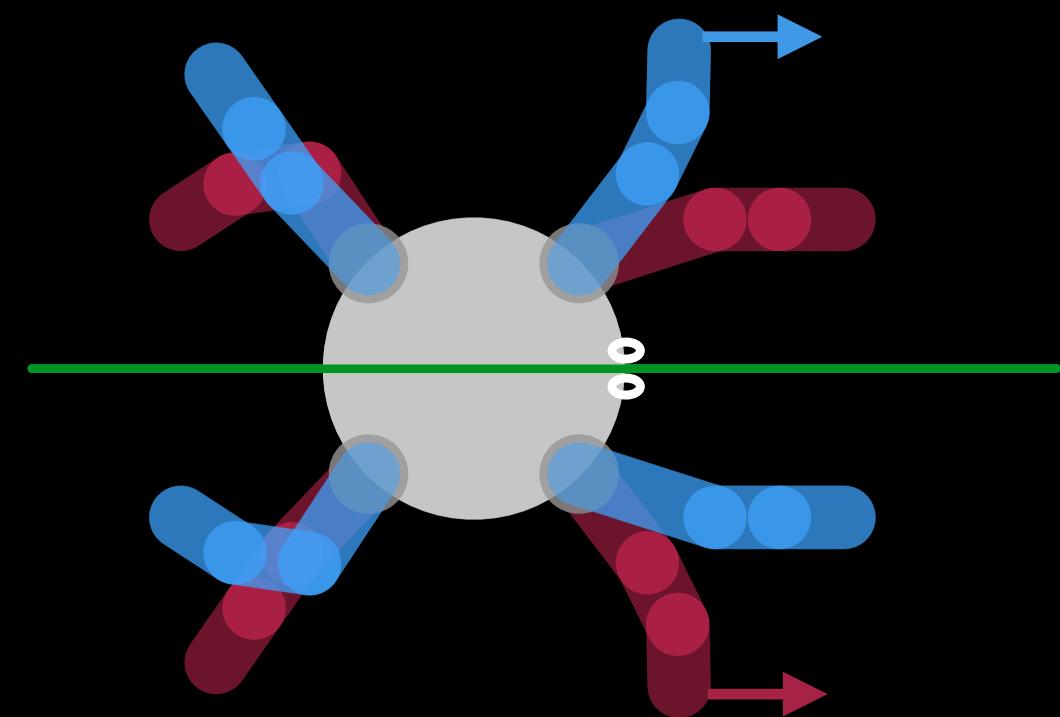
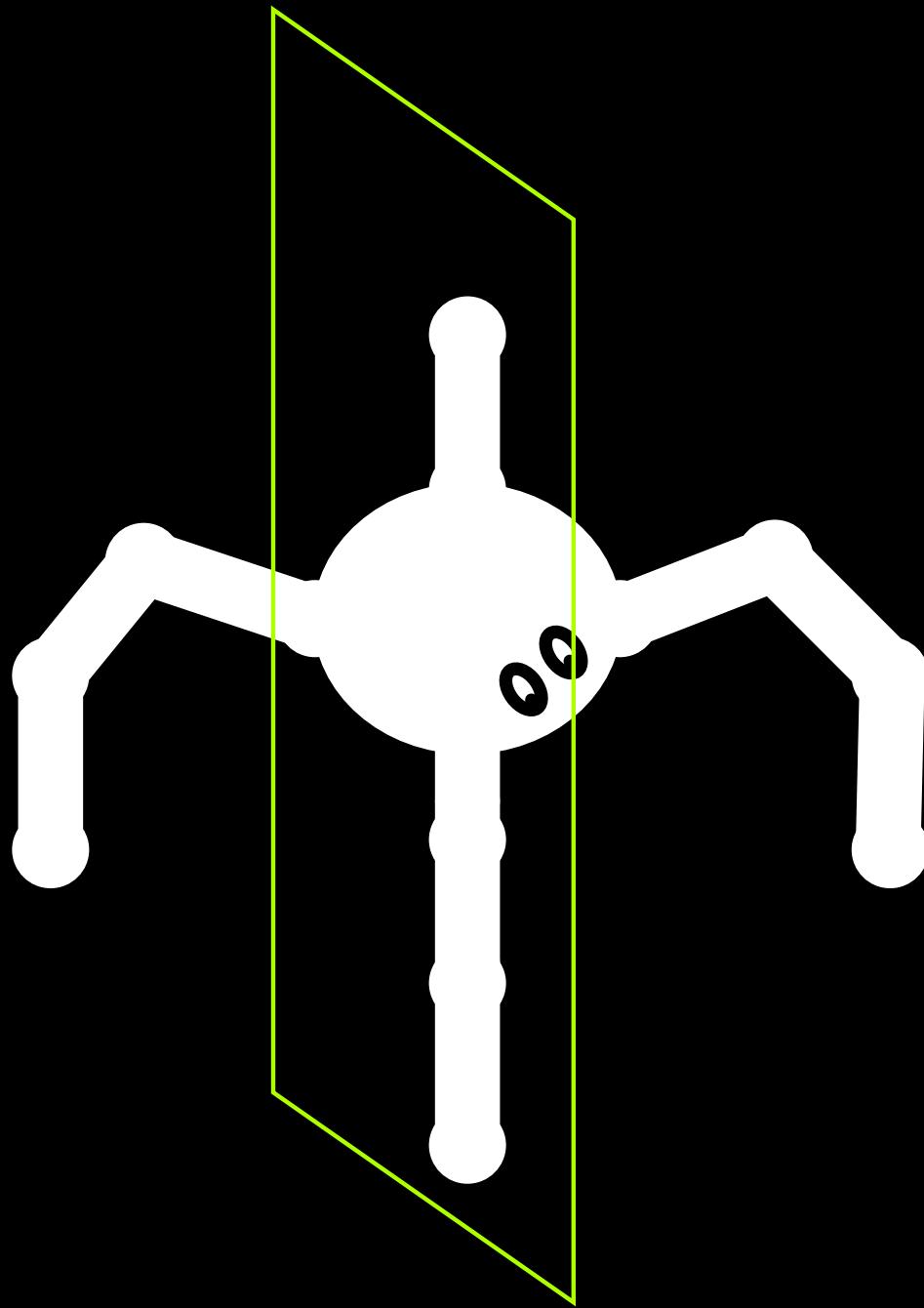


```
reward = move_reward * upright_reward
```

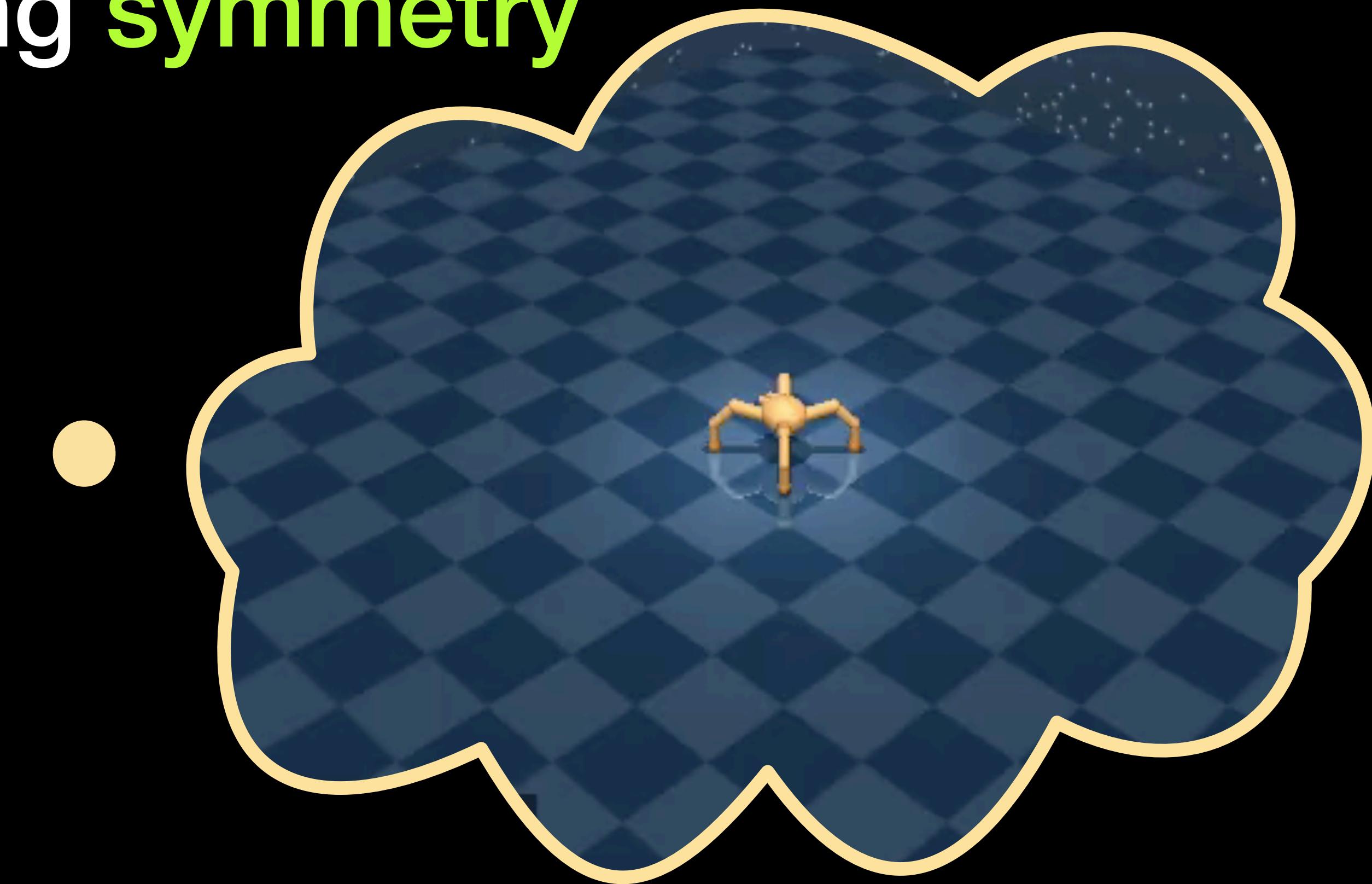
# Augmenting learning using symmetry in a biologically-inspired domain

# Symmetry in the Quadruped Domain

# Symmetry in the Quadruped Domain

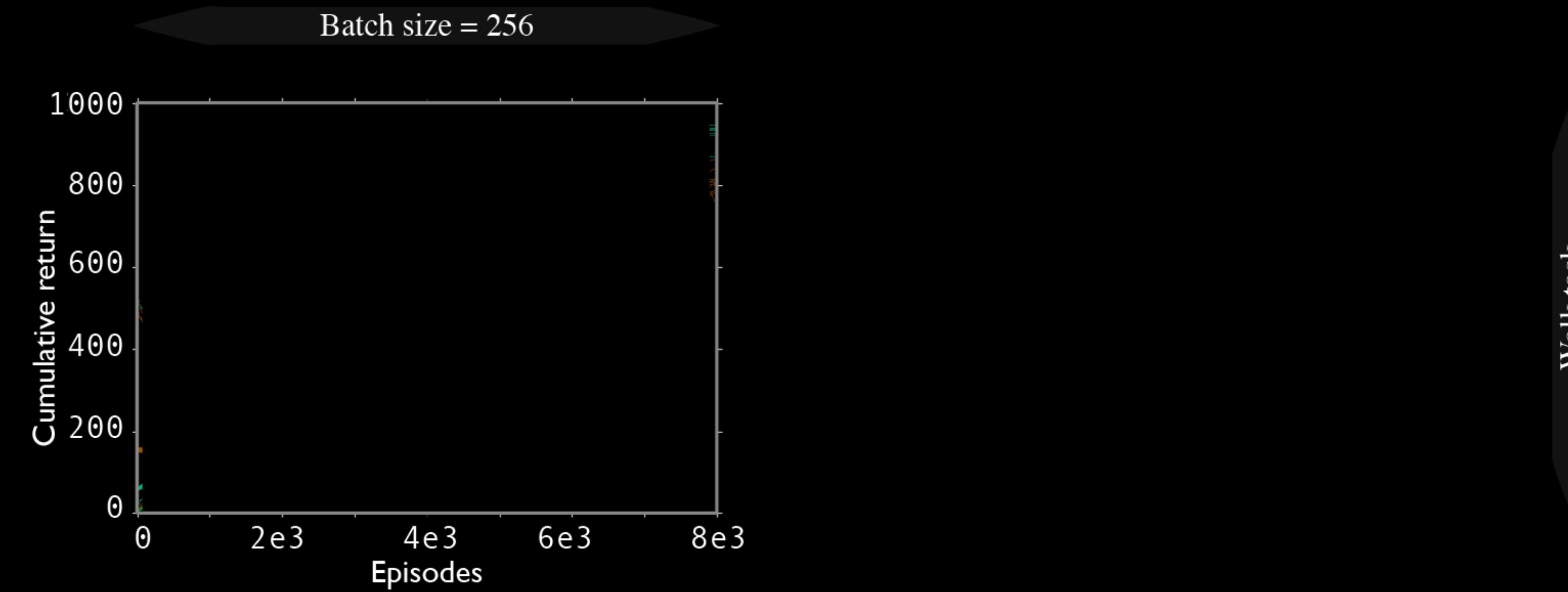


# Augmenting learning using symmetry



**Critic** augments the **batches** of data from the **actor** with **batches** of mirrored data.

# Preliminary results



# Preliminary results

Learning on data



Learning on **augmented** data



Agent that learns on **augmented** data looks more natural.

We need more quantitative analysis.

# Algorithm

---

```
1: given batch-size (K), number of actions (N), Q-function  $\hat{Q}$ , old-policy  $\pi_k$  and replay-buffer
2: initialize  $\pi_\theta$  from the parameters of  $\pi^{(k)}$ 
3: repeat
4:   Sample states with batch of size N from replay buffer
5:   Step 1: sample based policy (weights)
6:    $q(a_i|s_j) = q_{ij}$ , computed as:
7:   for  $j = 1, \dots, K$  do
8:     for  $i = 1, \dots, N$  do
9:        $a_i \sim \pi_k(a|s_j)$ 
10:       $Q_{ij} = \hat{Q}(s_j, a_i)$ 
11:       $q_{ij} \propto \exp(Q_{ij}/\text{temperature parameter}))$ 
12:    end for
13:  end for
14:  Calculate mirrored states and mirrored actions
15:  Step 2: update parametric policy
16:  Given the data-set  $\{s_j, (a_i, q_{ij})_{i=1 \dots N}\}_{j=1 \dots K}$  and  $\{s_j^{\text{mirrored}}, (a_i^{\text{mirrored}}, q_{ij})_{i=1 \dots N}\}_{j=1 \dots K}$ 
17:  Update the Policy by taking gradient of following weighted maximum likelihood objective
18:   $\max_\theta (\sum_j^K \sum_i^N q_{ij} \log \pi_\theta(a_i|s_j) + \sum_j^K \sum_i^N q_{ij} \log \pi_\theta(a_i^{\text{mirrored}}|s_j^{\text{mirrored}}))$ 
19:  (subject to additional (KL) regularization)
20: until Fixed number of steps
21: return  $\pi_{(k+1)}$ 
```

---

# Preliminary results

