

# **Efficient Meta-RL from Pixels**

**Kelly Marshall and Chingun Khasar 2021**

# High Level overview

Kelly & Chingun Fall 2021

# Recap: Meta- Reinforcement Learning



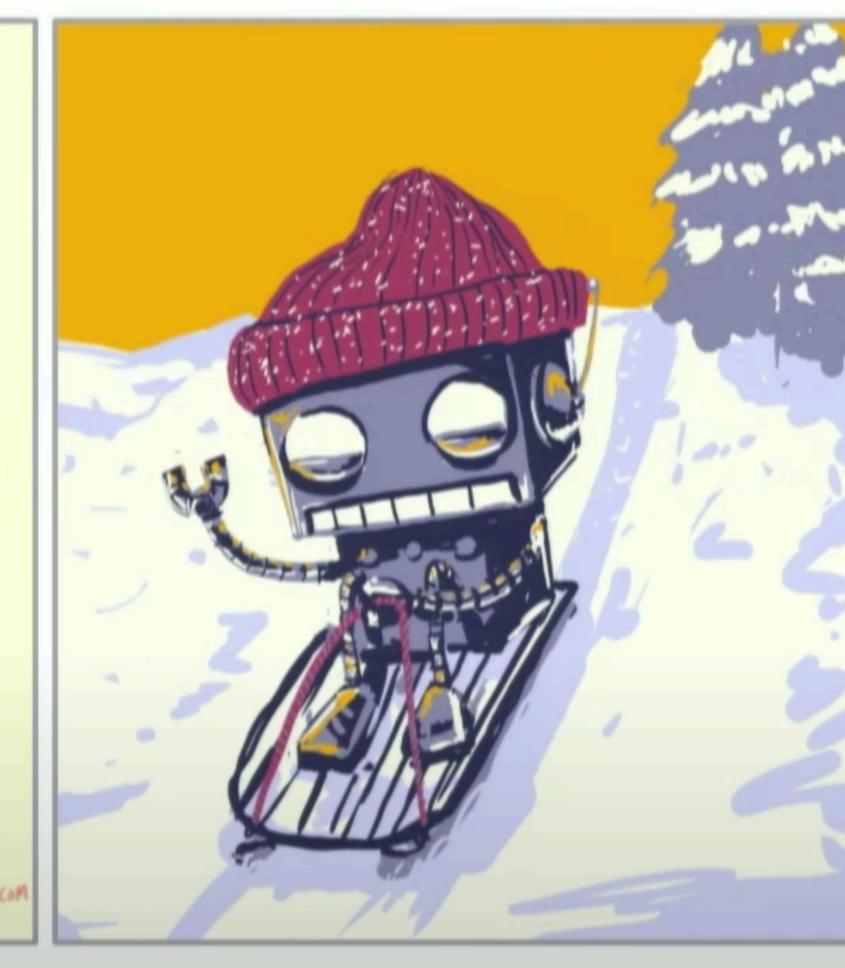
# Recap: Meta- Reinforcement Learning



M1



M2



M1



Mtest

# Recap: Meta- Reinforcement Learning

$$L(\pi_\theta) = \mathbb{E}_{\tau_i \sim p(\tau)} [\mathbb{E}_{s_t, a_t \sim \pi_{\theta_i}} [\sum_{t=1}^H \gamma^t r_t]]$$

# Meta-RL structure PEARL

## Probabilistic Embeddings for Actor-Critic Learning

---

**Algorithm 1** PEARL Meta Training
 

---

**Require:** Batch of training tasks  $\{\tau\}_{i=1..T}$  from  $p(\tau)$ , learning rates  $\alpha_1, \alpha_2, \alpha_3$

```

1: Initialize replay buffer  $\mathcal{B}^i$  for each training task
2: while not done do
3:   for each  $\tau_i$  do
4:     Initialize context  $\mathbf{c}^i = \{\}$ 
5:     for  $k = 1..K$  do
6:       Sample  $z \sim q_\phi(z|c^i)$ 
7:       Gather data from  $\pi_\theta(a|s, z)$  and add to  $\mathcal{B}_i$ 
8:       Update  $c_i = \{(s_j, a_j, r_j)\}_{j=1..N}$ 
9:     end for
10:   end for
11:   for each step in training steps do
12:     for each  $\tau_i$  do
13:       Sample context  $c^i \sim \mathcal{S}_c(\mathcal{B}_i)$  and RL batch  $b^i \sim \mathcal{B}_i$ 
14:       Sample  $z \sim q_\phi(z|c^i)$ 
15:        $\mathcal{L}_{actor}^i = \mathcal{L}_{actor}(b^i, \mathbf{z})$ 
16:        $\mathcal{L}_{critic}^i = \mathcal{L}_{critic}(b^i, \mathbf{z})$ 
17:        $\mathcal{L}_{KL}^i = \beta D_{KL}(q_\phi(\mathbf{z}|\mathbf{c}^i) || r(\mathbf{z}))$ 
18:     end for
19:      $\phi \leftarrow \phi - \alpha_1 \nabla_\phi \sum_i (\mathcal{L}_{critic}^i + \mathcal{L}_{KL}^i)$ 
20:      $\theta_\pi \leftarrow \theta_\pi - \alpha_2 \nabla_\theta \sum_i \mathcal{L}_{actor}^i$ 
21:      $\theta_Q \leftarrow \theta_Q - \alpha_3 \nabla_\theta \sum_i \mathcal{L}_{critic}^i$ 
22:   end for
23: end while
  
```

---

**Algorithm 2** PEARL Meta Testing
 

---

**Require:** Test task  $\tau \sim p(\tau)$

```

1: Initialize context  $\mathbf{c}^\tau = \{\}$ 
2: for  $k=1..K$  do
3:   Sample  $z \sim q_\phi(\mathbf{z}|c^\tau)$ 
4:   Roll out policy  $\pi_\theta(\mathbf{a}|\mathbf{s}, \mathbf{z})$  to collect data  $D_k^\tau = \{\mathbf{s}_j, \mathbf{a}_j, \mathbf{s}', r_j\}$ 
5:   Accumulate context  $\mathbf{c}^\tau = \mathbf{c}^\tau \cup D_k^\tau$ 
6: end for
  
```

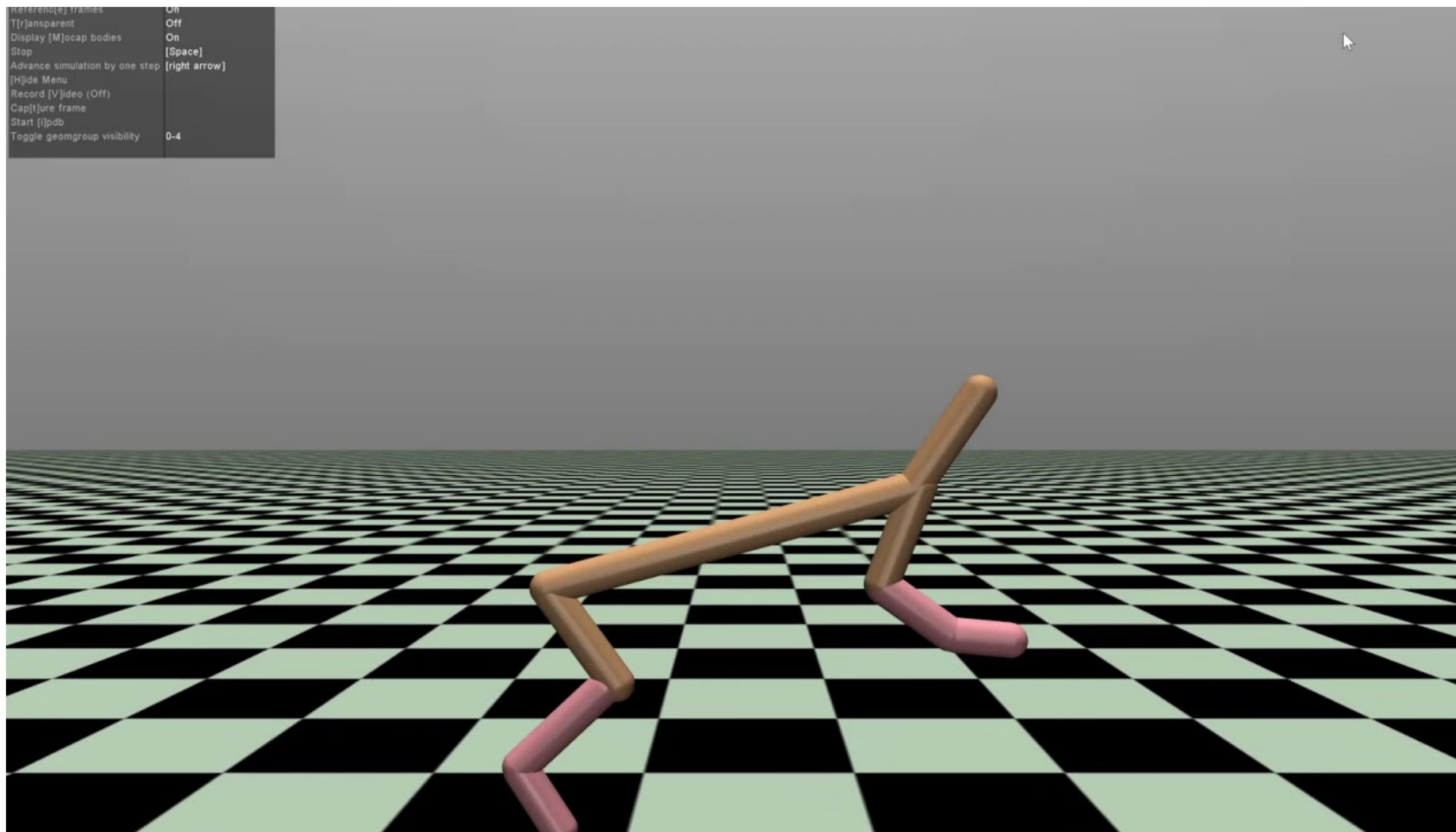
---

# Our Additions

Kelly & Chingun Fall 2021

# Image + meta-RL

## On Half-Cheetah



# Design Questions of Image + meta-RL

- What loss function to use?
- How to balance the use of true state representations?

# CNN + PEARL

1. Train CNN model image to state classification.
2. Train RL using state vectors.
3. During test use CNN predicted vectors

# Distribution Shift from +Image

CNN error + RL error = bad Agent

# CNN + PEARL

1. Train CNN model for understanding image to state classification.
2. Train RL ~~using state vectors~~ **CNN predicted vectors.**
3. During test use CNN predicted vectors.

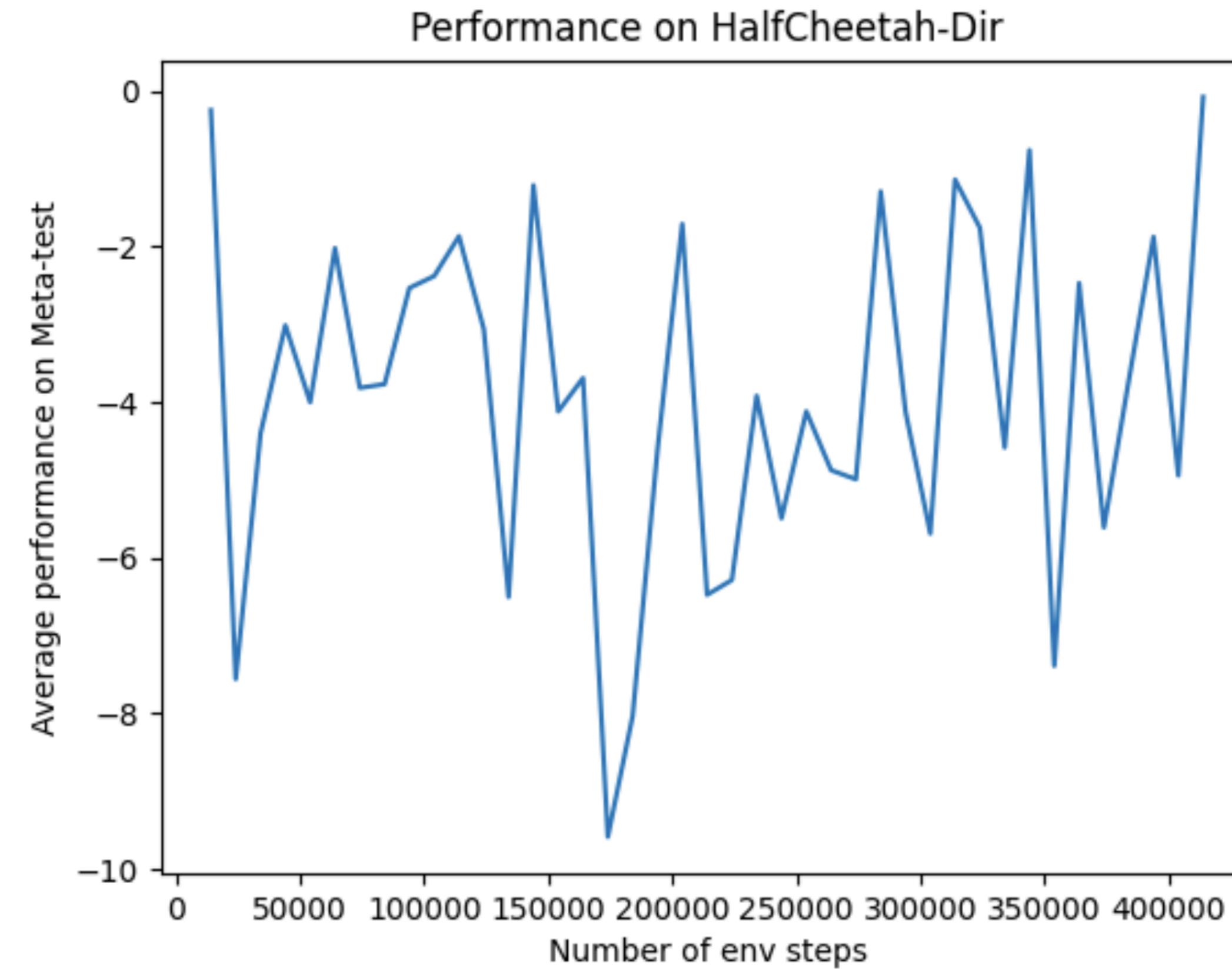
# **Augment CNN loss**

CNN loss = image classification loss + RL reward loss

# CNN + PEARL

1. Train CNN model ~~for understanding image to state classification~~ with **added signal from RL reward.**
2. Train RL ~~using state vectors~~ **CNN predicted vectors.**
3. During test use CNN predicted vectors.

# Experiments



**Thank you**