

2024951308 Robin Henri Dieudonné

DATA403(00) Student Presentation
(Reinforcement Learning)

Contents

1. Selected environments
2. Selected Algorithm(s)
3. Results
4. Conclusion

1. Selected Environments

- 2 Environments

- Ant (Lv.2)
- HumanoidStandup (Lv.3)

- I selected them because it is like watching children grow up, which is fun

2. Selected Algorithm(s)

□ PPO for Ant-v4 & HumanoidStandup-v4

Algorithm 1 PPO-Clip

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 4: Compute rewards-to-go \hat{R}_t .
- 5: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 6: Update the policy by maximizing the PPO-Clip objective:

$$\theta_{k+1} = \arg \max_{\theta} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \min \left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)} A^{\pi_{\theta_k}}(s_t, a_t), \quad g(\epsilon, A^{\pi_{\theta_k}}(s_t, a_t)) \right),$$

typically via stochastic gradient ascent with Adam.

- 7: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 8: **end for**
-

2. Selected Algorithm(s)

□ PPO hyperparameters for Ant-v4 & HumanoidStandup-v4

<i>batch_size=64</i>	<i>vf_coeff=0.5</i>
<i>n_steps=2048</i>	<i>max_grad_norm=0.5</i>
<i>lr=3e-4 (default)</i>	<i>gae_lambda=0.95</i>
<i>buffer_size=1e6</i>	<i>gamma=0.99</i>
<i>entropy_coeff=0.0 (default)</i>	<i>clip_range=0.2</i>

We use a small learning rate with the Adam optimizer as we train for several steps
Parameters that we play with are the learning rate and the entropy coefficient.

2. Selected Algorithm(s)

□ TD3/Twin-Delayed DDPG for Ant-v4

Algorithm 1 Twin Delayed DDPG

```

1: Input: initial policy parameters  $\theta$ , Q-function parameters  $\phi_1, \phi_2$ , empty replay buffer  $\mathcal{D}$ 
2: Set target parameters equal to main parameters  $\theta_{\text{targ}} \leftarrow \theta, \phi_{\text{targ},1} \leftarrow \phi_1, \phi_{\text{targ},2} \leftarrow \phi_2$ 
3: repeat
4:   Observe state  $s$  and select action  $a = \text{clip}(\mu_\theta(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$ , where  $\epsilon \sim \mathcal{N}$ 
5:   Execute  $a$  in the environment
6:   Observe next state  $s'$ , reward  $r$ , and done signal  $d$  to indicate whether  $s'$  is terminal
7:   Store  $(s, a, r, s', d)$  in replay buffer  $\mathcal{D}$ 
8:   If  $s'$  is terminal, reset environment state.
9:   if it's time to update then
10:     for  $j$  in range(however many updates) do
11:       Randomly sample a batch of transitions,  $B = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$ 
12:       Compute target actions
          
$$a'(s') = \text{clip}(\mu_{\theta_{\text{targ}}}(s') + \text{clip}(\epsilon, -c, c), a_{\text{Low}}, a_{\text{High}}), \quad \epsilon \sim \mathcal{N}(0, \sigma)$$

13:       Compute targets
          
$$y(r, s', d) = r + \gamma(1 - d) \min_{i=1,2} Q_{\phi_{\text{targ},i}}(s', a'(s'))$$


```

Update Q-functions by one step of gradient descent using

$$\nabla_{\phi_i} \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_{\phi_i}(s, a) - y(r, s', d))^2 \quad \text{for } i = 1, 2$$

if $j \bmod \text{policy_delay} = 0$ **then**
 Update policy by one step of gradient ascent using

$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi_1}(s, \mu_\theta(s))$$

Update target networks with

$$\begin{aligned} \phi_{\text{targ},i} &\leftarrow \rho \phi_{\text{targ},i} + (1 - \rho) \phi_i \\ \theta_{\text{targ}} &\leftarrow \rho \theta_{\text{targ}} + (1 - \rho) \theta \end{aligned} \quad \text{for } i = 1, 2$$

```

18:   end if
19: end for
20: end if
21: until convergence

```

3. Results

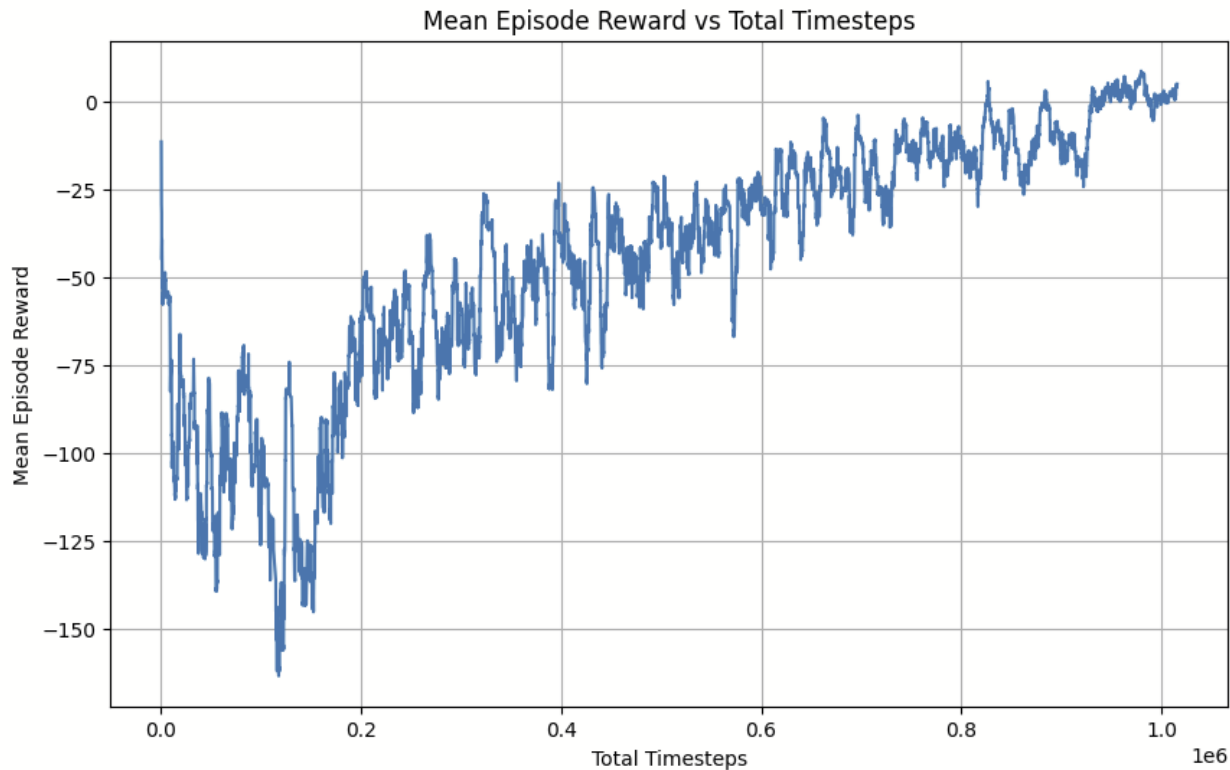
□ Ant-v4

Proximal Policy Optimization (PPO)

ENTROPY = 0.0

LEARNING RATE = $3e-4$

FINAL EPISODE MEAN REWARD = 4.69



3. Results

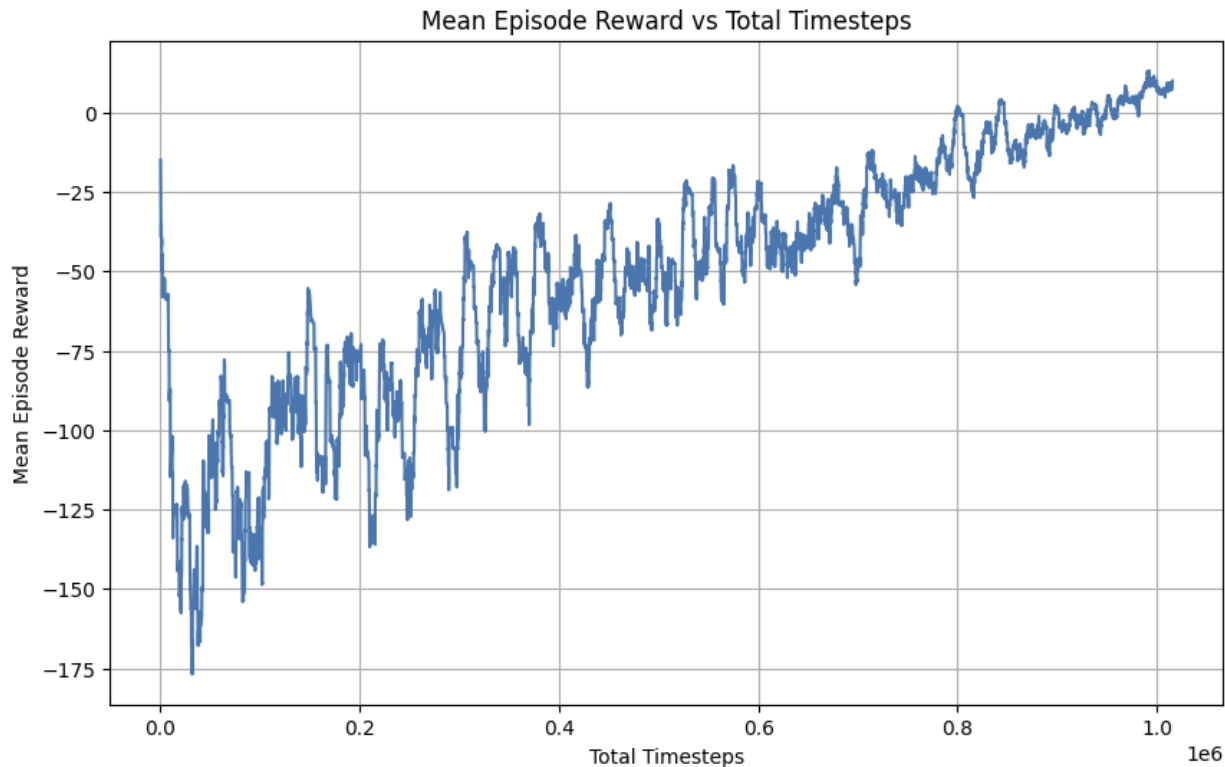
□ Ant-v4

Proximal Policy Optimization (PPO)

ENTROPY = 0.001

LEARNING RATE = $3e-4$

FINAL EPISODE MEAN REWARD = 11.9



3. Results

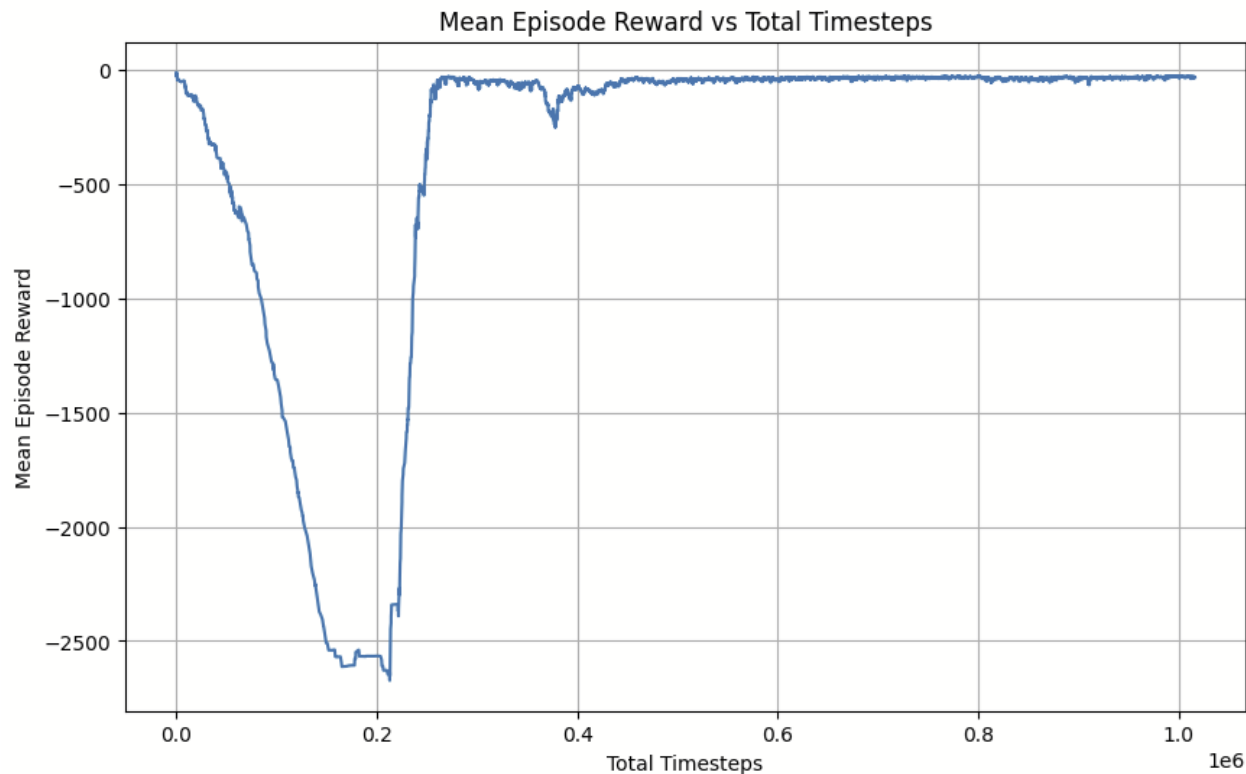
□ Ant-v4

Proximal Policy Optimization (PPO)

ENTROPY = 0.0

LEARNING RATE = $9\text{e-}3$

FINAL EPISODE MEAN REWARD = -34.1



3. Results

□ Ant-v4

TD3/Twin-Delayed DDPG for Ant-v4

learning_starts=10000

batch_size=100

learning_rate=1e-3

gamma=0.99

tau=0.005

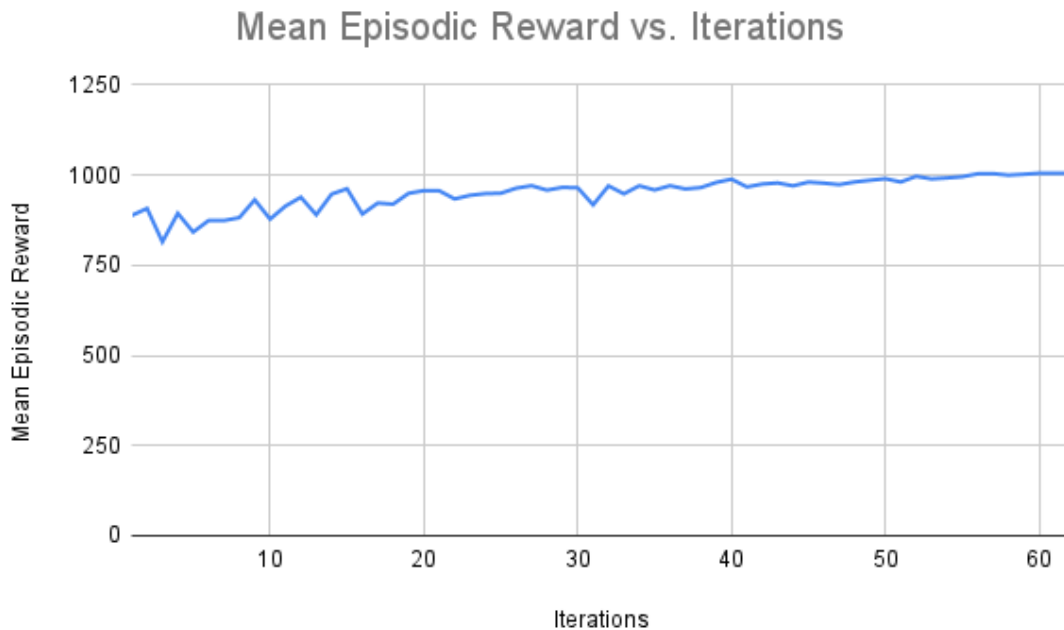
gradient_steps=-1

policy_delay=2

target_policy_noise=0.2

target_noise_clip=0.5

mean reward = 1.004e+03



3. Results

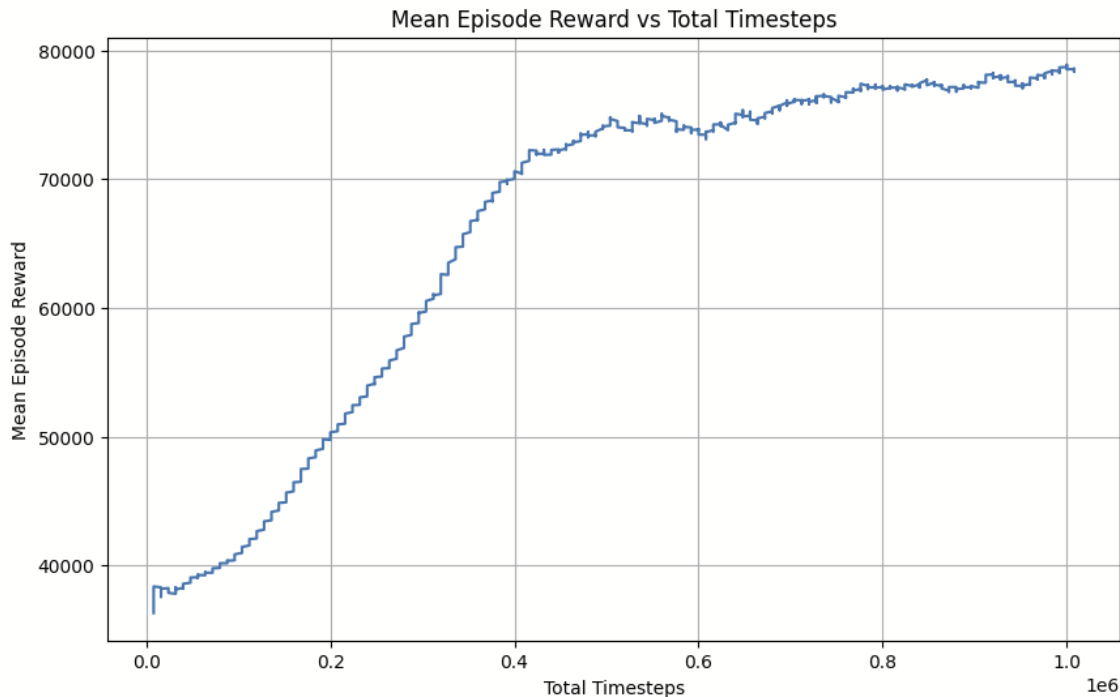
□ Humanoid Standup-v4

Proximal Policy Optimization (PPO)

ENTROPY = 0.0

LEARNING RATE = $3e-4$

FINAL EPISODE MEAN REWARD = $7.84e+04$



3. Results

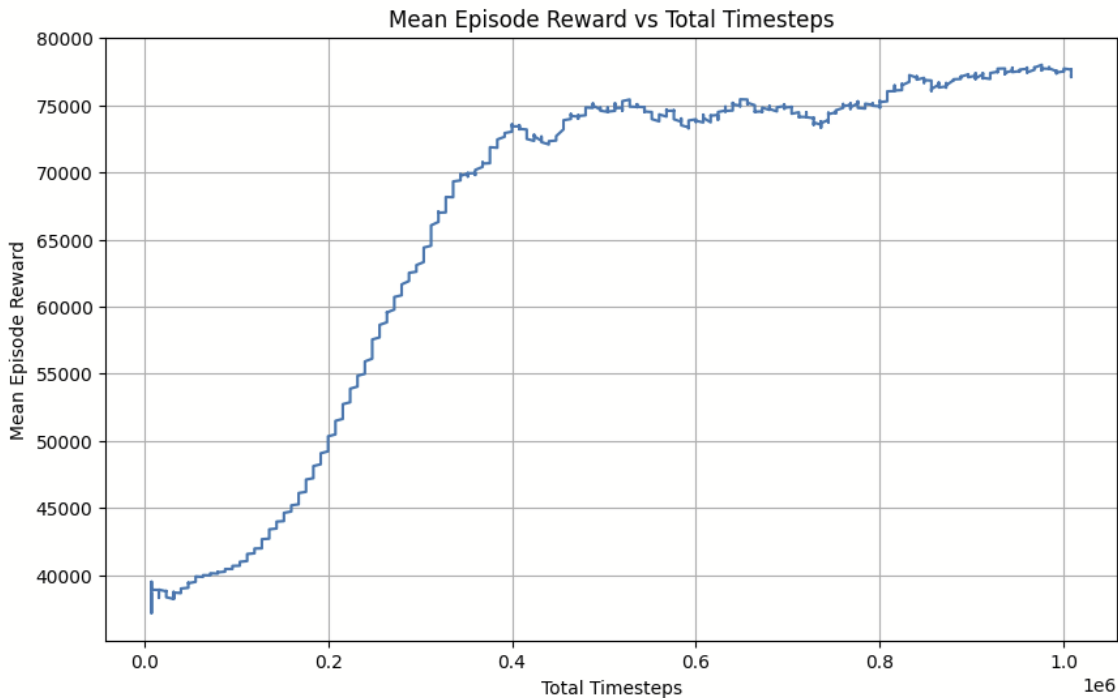
□ Humanoid Standup-v4

Proximal Policy Optimization (PPO)

ENTROPY = 0.001

LEARNING RATE = $3e-4$

FINAL EPISODE MEAN REWARD = $7.7e+4$



3. Results

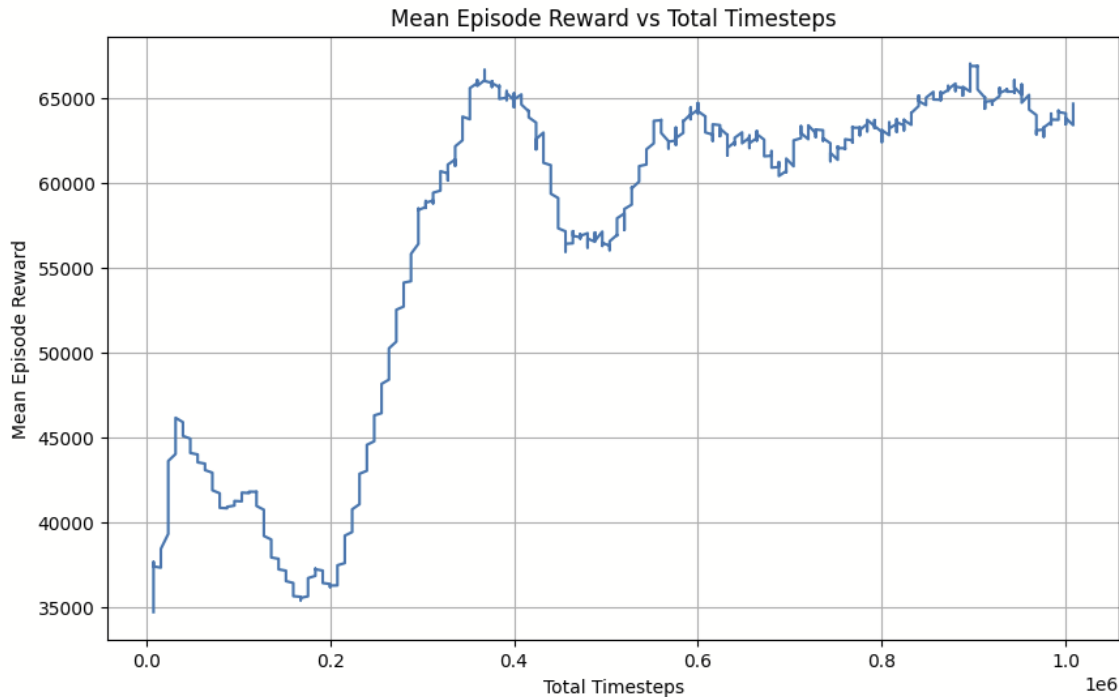
□ Humanoid Standup-v4

Proximal Policy Optimization (PPO)

ENTROPY = 0.0

LEARNING RATE = $9e-3$

FINAL EPISODE MEAN REWARD = $6.48e+04$



Conclusion

- ❑ Results are very sensitive to hyperparameters choices
- ❑ Well-balanced entropy coeff => better results
- ❑ PPO approx. 4x faster than TD3
- ❑ TD3 better than PPO for Ant-V4

Thanks!

- Q & A