

Leveraging Proximal Policy Optimization for Efficient Robotic Manipulation

ECE 517: Reinforcement Learning

Group 7

Shovan Chowdhury, Sk Hasibul Alam, Umarbek Guvercin, Babita Babita

Main Goal:

- Train a robotic arm for reaching and lifting tasks using reinforcement learning (RL).

Specific Objectives:

- Develop a stable and efficient policy for robotic control.
- Use Proximal Policy Optimization (PPO) to handle high-dimensional continuous action spaces.

Current State of the Work:

- Robotics control challenges: high-dimensional action spaces, continuous control, and non-linear dynamics.
- Traditional approaches: Supervised learning and rule-based systems.
- Emergence of RL: Shift to learning policies directly from interactions with the environment.

Why Reinforcement Learning?

- Self-learning and adaptability.
- No need for explicit programming of all behaviors.

Progress in RL for Robotics:

- DDPG, SAC, and TD3 for continuous control.
- Limitations: unstable training, exploration inefficiencies.

PPO:

- Stability through clipped updates.
- Effective for complex robotic tasks.

Why PPO?

COMPARISON OF PPO WITH OTHER RL ALGORITHMS

Benchmark	DQN	A3C/A2C	TRPO	PPO
Handles continuous actions	✗	✓	✓	✓
Stability in training	✗	✗	✓	✓
Computational efficiency	✓	✓	✗	✓
Simplicity of implementation	✓	✓	✗	✓
Robustness to hyperparameters	✗	✗	✓	✓
Suitable for robotics	✗	✓	✓	✓

PPO Algorithm

During PPO learning, $GAE(\lambda)$ improves the estimation of advantage function.

- Initialize advantage estimate A for all timesteps
- Calculate TD error for each timestep
- Iteratively calculate A for each timestep
- Return the normalized A

λ controls the balance between bias and variance.

Algorithm 1 Generalized Advantage Estimation (GAE)

```
1: function GAE( $r, d, V, V'_{last}$ )
2:    $a \leftarrow 0$ 
3:    $A \leftarrow \text{zeros}(r)$ 
4:   for reverse  $i \leftarrow$  number of rows in  $r$  do
5:     if  $i \neq$  last row of  $r$  then
6:        $V'_i \leftarrow V_{i+1}$ 
7:     else
8:        $V'_i \leftarrow V'_{last}$ 
9:     end if
10:     $a \leftarrow r_i - V_i + \gamma (\neg d_i) (V'_i - \lambda a)$ 
11:     $A_i \leftarrow a$ 
12:  end for
13:  return  $R \leftarrow A + V$ 
14:  return  $A \leftarrow (A - \bar{A}) / (A_\sigma + 10^{-8})$ 
15: end function
```

Policy Update:

For each iteration:

- Collect, in a rollout memory, a set of states s , actions a , rewards r , dones d , log probabilities $\log p$ and values V on policy using π_θ and V_ϕ
- Estimate returns R and advantages A using $\text{GAE}(\lambda)$ from the collected data $[r, d, V]$
- Compute the clipped surrogate objective (policy loss) with $ratio$ as the probability ratio between the action under the current policy and the action under the previous policy:

$$L_{\pi_\theta}^{clip} = \mathbb{E} \left[\min(A \cdot ratio, A \cdot \text{clip}(ratio, 1 - \epsilon, 1 + \epsilon)) \right]$$
- Compute the value loss L_{V_ϕ} as the MSE between the predicted values $V_{predicted}$ and the estimated returns R
- Optimize the total loss $L = L_{\pi_\theta}^{clip} - c_1 L_{V_\phi}$

Algorithm 2 Update for learning

```

1:  $V'_{last} \leftarrow V_\phi(s')$ 
2:  $R, A \leftarrow \text{GAE}(r, d, V, V'_{last})$ 
3:  $[[s, a, \log p, V, R, A]] \leftarrow \text{states, actions, log\_prob, values, returns, advantages}$   $\triangleright$  sample mini-batches
4: for all learning epochs do
5:   for all mini batches do
6:      $\log p' \leftarrow \pi_\theta(s, a)$ 
7:      $ratio \leftarrow \log p' - \log p$ 
8:      $KL_{divergence} \leftarrow \frac{1}{N} \sum_{i=1}^N (e^{ratio} - 1 - ratio)$ 
9:     if  $KL_{divergence} > KL_{threshold}$  then
10:       break
11:     end if
12:      $ratio \leftarrow e^{\log p' - \log p}$ 
13:      $L_{surrogate} \leftarrow A \cdot ratio$ 
14:      $L_{clipped} \leftarrow A \cdot \text{clip}(ratio, 1 - \epsilon, 1 + \epsilon)$ 
15:      $L_{\pi_\theta}^{clip} \leftarrow -\frac{1}{N} \sum_{i=1}^N \min(L_{surrogate}, L_{clipped})$ 
16:      $V_{predicted} \leftarrow V_\phi(s)$ 
17:      $V_{predicted} \leftarrow V + \text{clip}(V_{predicted} - V, -\beta, \beta)$ 
18:      $L_{V_\phi} \leftarrow \frac{1}{N} \sum_{i=1}^N (R - V_{predicted})^2$ 
19:     reset OPTIMIZER $_{\theta, \phi}$   $\triangleright$  torch Adam
20:      $\nabla_{\theta, \phi} (L_{\pi_\theta}^{clip} + L_{V_\phi})$ 
21:     clip_grad_norm( $\|\nabla_{\theta, \phi}\|$ )
22:     step OPTIMIZER $_{\theta, \phi}$ 
23:   end for
24: end for

```

➤ Environments

- Isaac-Reach-Franka-v0
- Isaac-Lift-Cube-Franka-v0

➤ Memories

- Rollout buffers → Random memories from SKRL

➤ Models

- Agent's brain → Neural networks as function approximator

➤ Reward

Reach		Lift	
Name	Weight	Name	Weight
end_effector_pos_track	-0.2	reaching_object	+1
end_effector_pos_track_fine	+0.1	lifting_object	+15
end_effector_orientation_track	-0.1	obj_goal_track	+16
action_rate	-0.0001	obj_goal_track_fine	+5
joint_vel	-0.0001	action_rate	-0.0001
		joint_vel	-0.0001

➤ Training [50 parallel instances]

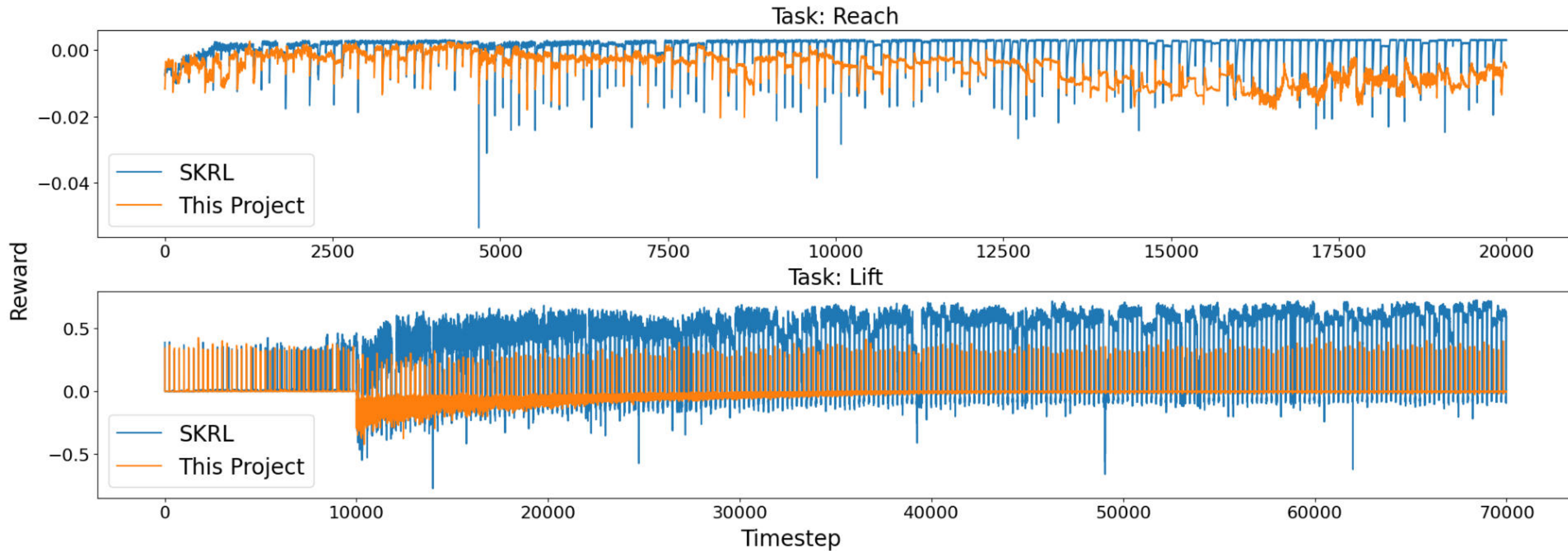
- 20000 timesteps for Reach
- 70000 timesteps for Lift

➤ Evaluation [1 instance]

- 2000 timesteps → Video rendered

Results

(1) Learning Rate Scheduler *and* Preprocessor

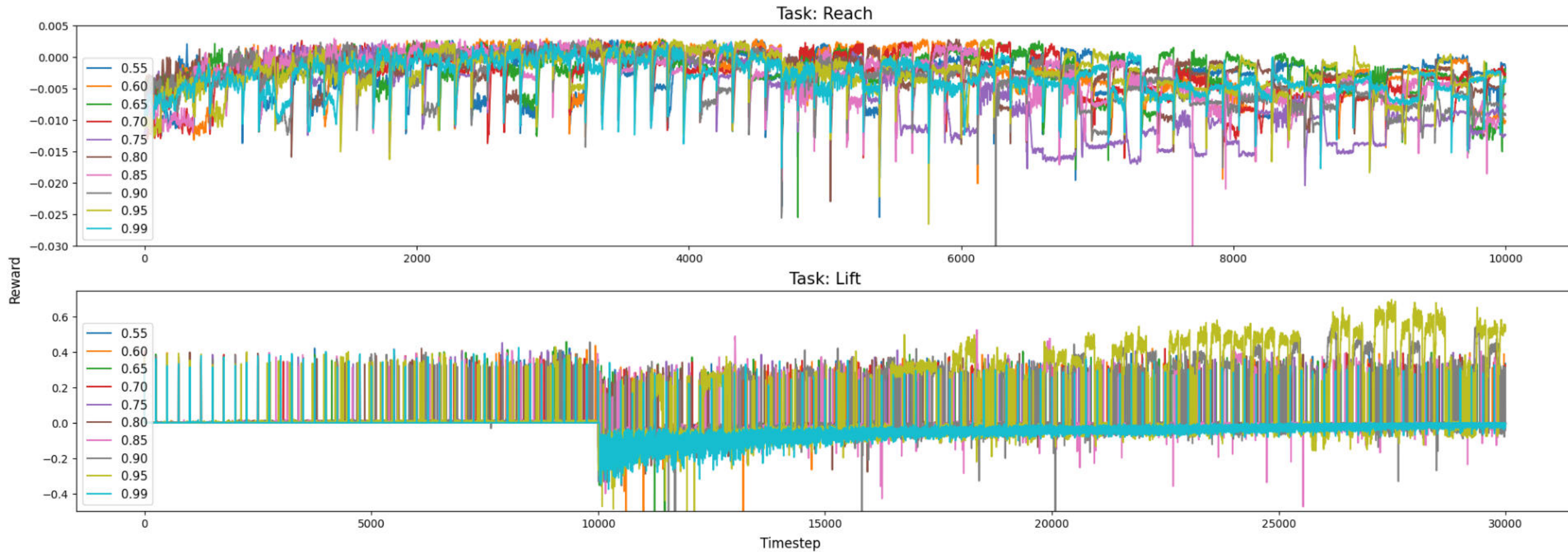


Takeaway:

- High η helps earlier in training (exploration). Low η helps later in training (exploitation).
- NN used in PPO may work better when inputs are standardized.

Results

(2) Discount Factor

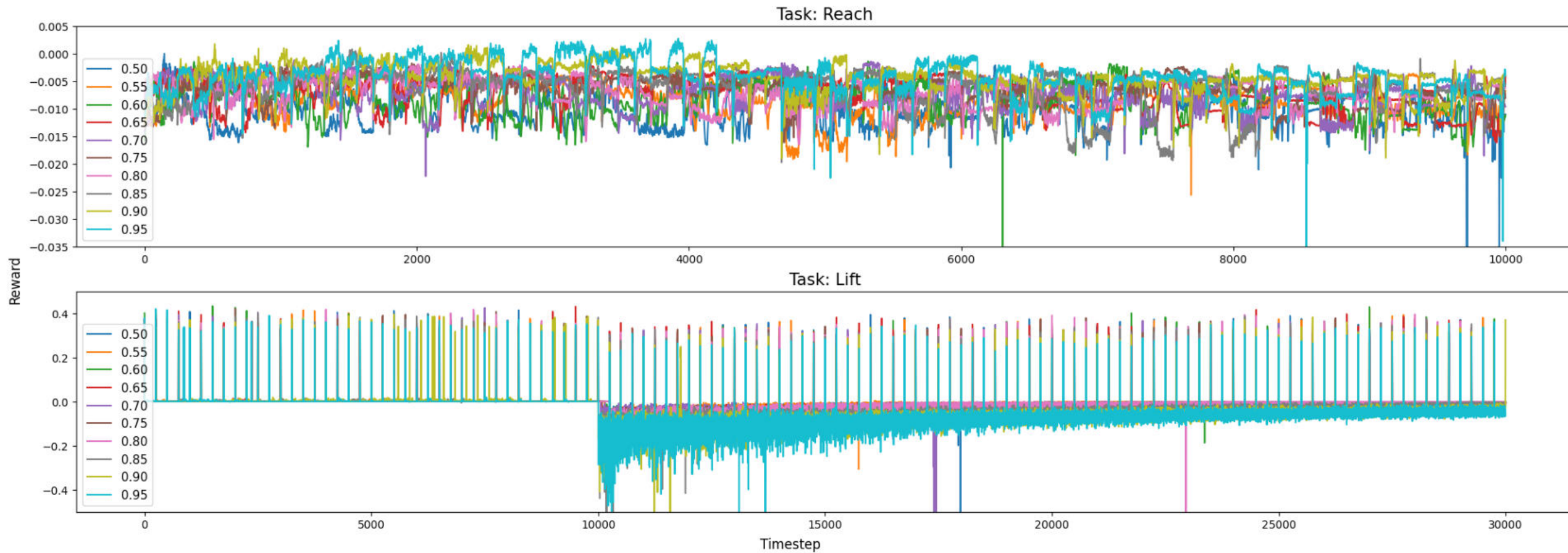


Takeaway:

- High γ considers long-term rewards (exploration), excels in sparse reward env. [\rightarrow helped us]
- Low γ considers immediate rewards (exploitation), struggles in sparse reward env.

Results

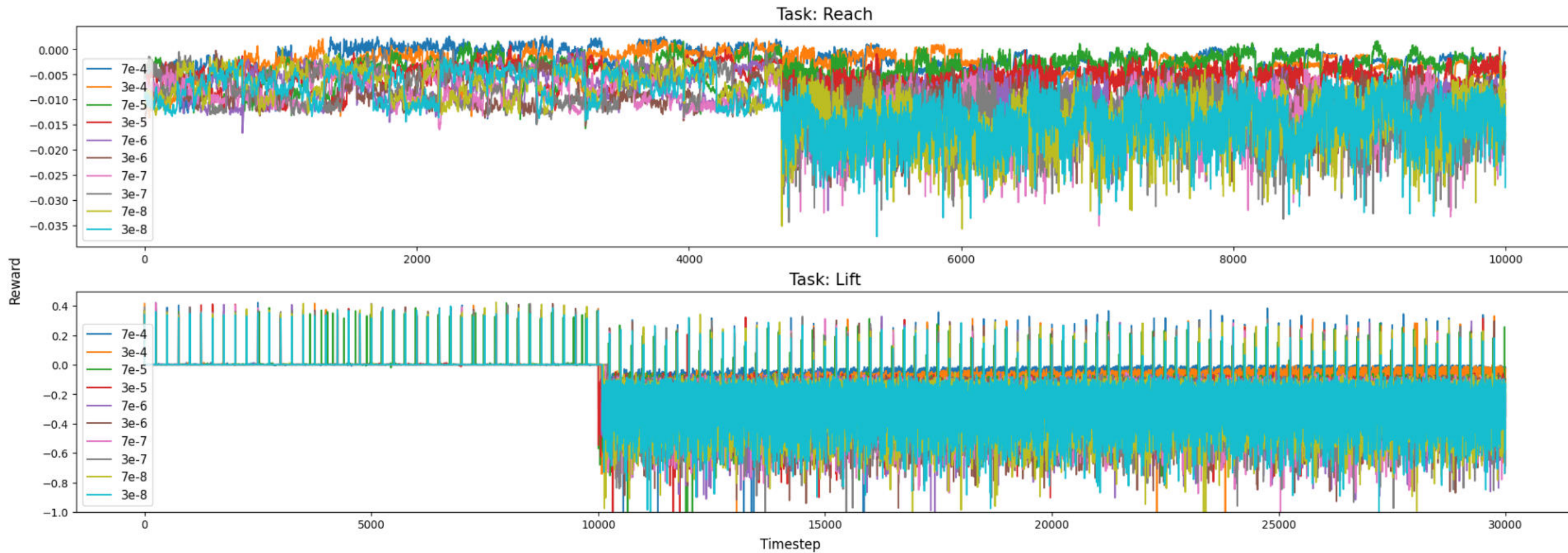
(3) TD Coefficient



Takeaway:

- High λ considers long-term rewards (exploration), reduces bias. [→ helped us]
- Low λ considers immediate rewards (exploitation), reduces variance.

(4) Learning Rate

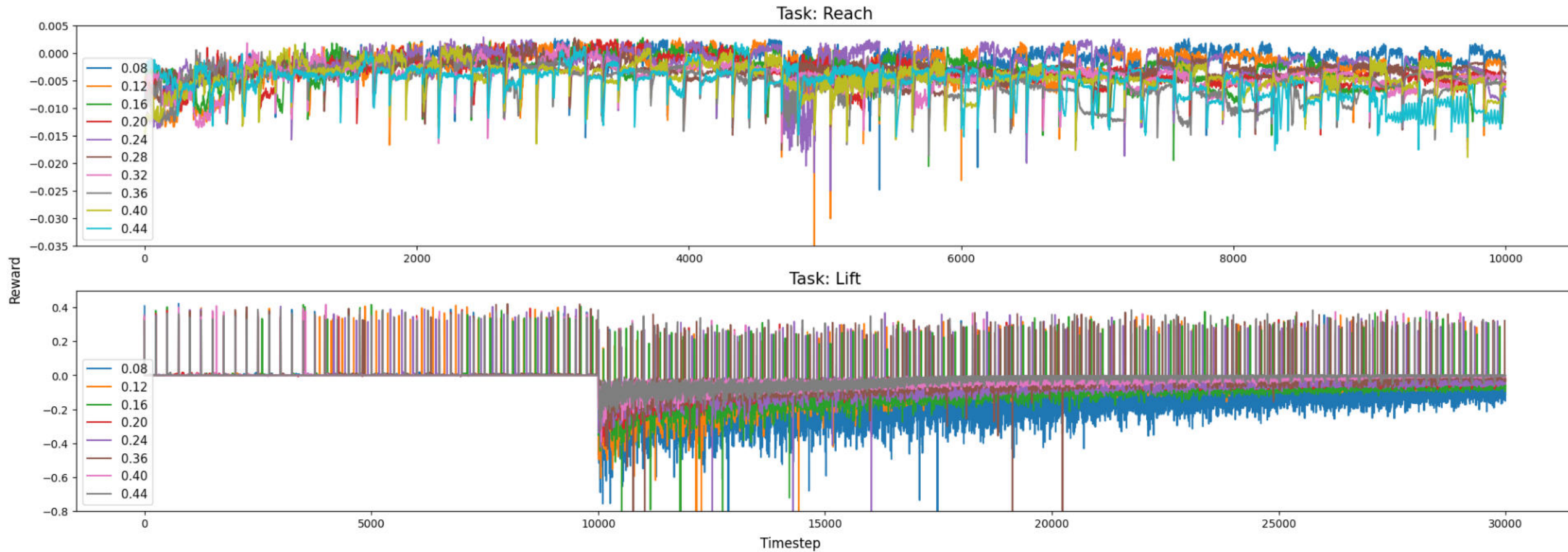


Takeaway:

- High η promotes rapid policy change (exploration), converges faster. [\rightarrow helped us]
- Low η promotes gradual policy refinement (exploitation), requires more timesteps.

Results

(5) Clip Ratio



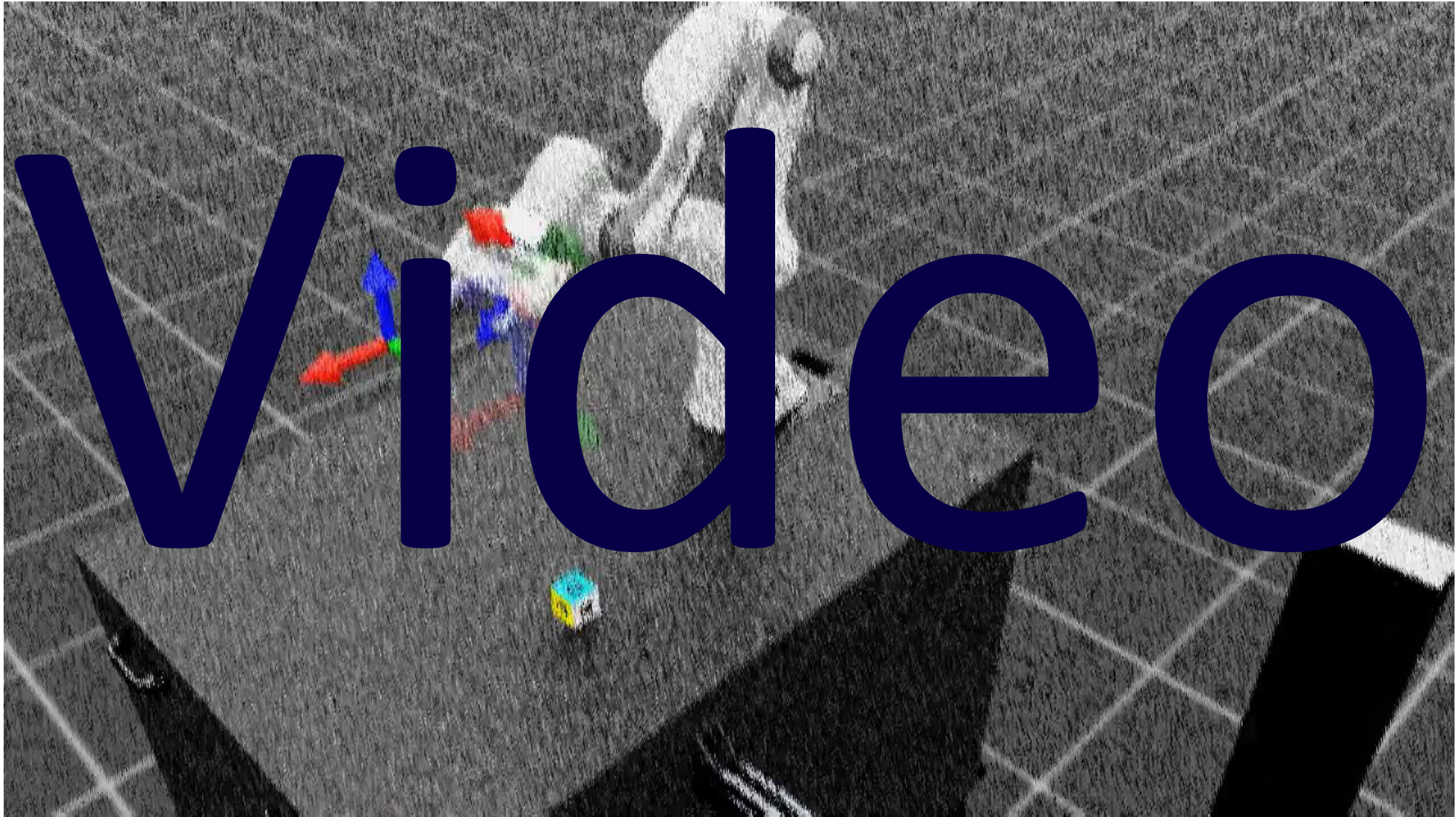
Takeaway:

- High ϵ allows larger policy update (exploration), risks overshooting.
- Low ϵ prevents large change (exploitation), mitigates noisy estimates. [→ helped us]

(1) Reach



(2) Lift



- The project was a success!
 - η -scheduler and preprocessor might have helped.
- Goal: Explore various hyperparameters
 - not reinvent the PPO algorithm itself
- Higher values of γ , λ , η helped.
 - promoted exploration, but increased variance
- Lower value of ϵ helped.
 - promoted exploitation, but required more iterations

- Shovan – memory and models, literature review, writing report and preparing slides
- Hasibul – interfacing modules, exploring hyperparameters, rendering video, writing report, preparing slides
- Umarbek – PPO update function, report
- Babita – GAE and PPO update function, report

Thank you!

Questions? Comments? Concerns?