

A ALGORITHM

Algorithm 1 FLR with HER

```

1: Sample the batch of transition  $\{(s_t, a_t, s_{t+1}, sg_i)\}_{|B|}$  in replay buffer  $\mathcal{B}$ 
2: Sample an additional goals  $\phi(s_{t+F})$  from current episode with a probability of  $p$ 
3:  $sg_i \leftarrow \phi(s_{t+F})$ 
4: if  $\|\phi(s_{t+1}) - sg_i\| < \delta$  then
5:    $r_t \leftarrow \max_{sg: \varepsilon \leq \|sg_i - sg\|} \alpha V_{\pi_l}^l(s_{t+1}, sg)$ 
6: else
7:    $r_t \leftarrow -1$ 
8: end if
9: Update  $Q_\theta$  using the transition  $(s_t, a_t, r_t, s_{t+1}, sg, sg')$ 

```

Algorithm 2 DHRL (Lee et al., 2022)

```

1: Input: initial random steps  $\tau_{\text{randomwalk}}$ , initial steps without planning  $\tau_{\text{w/o graph}}$ , total training
   step  $\tau_{\text{total}}$ , Env, low-level agent  $Q_{\text{critic}, \theta_1}^{\text{lo}}, Q_{\text{graph}, \theta_2}^{\text{lo}}$  and  $\pi_{\phi_1}^{\text{lo}}$ , high-level agent  $Q_{\theta_3}^{\text{hi}}$  and  $\pi_{\phi_2}^{\text{hi}}$ 
2:  $Dist(s, g) := \log_\gamma(1 + (1 - \gamma)Q_{\text{graph}, \theta_2}^{\text{lo}}(s, \pi(s, g)|g))$ 
3: for  $\tau = 1$  to  $\tau_{\text{total}}$  do
4:   if Env.done then
5:     Env.reset (episode step resets to 0)
6:   end if
7:   if  $\tau < \tau_{\text{randomwalk}}$  then
8:      $a_t \leftarrow \text{random.uniform}(\text{high} = \text{action.high}, \text{low} = \text{action.low})$   $\triangleleft$  initial random
    rollout
9:   else if  $\tau < \tau_{\text{w/o graph}}$  then
10:     $a_t \leftarrow \text{vanilla HRL}(sg_t = \pi_{\phi_2}^{\text{hi}}(s_t, g) \text{ and } \pi_{\phi_1}^{\text{lo}}(s_t, sg_t))$   $\triangleleft$  act without planning
11:   else
12:     if Graph  $\mathbf{G}$  is not initialized then
13:       Create a graph  $\mathbf{G}(\mathbf{V}, \mathbf{E})$  using FPS algorithm  $\triangleleft$  initialize graph
14:     end if
15:     if episode step(the step of the environment) %  $c_l = 0$  then
16:        $sg_t \leftarrow \pi_{\phi_2}^{\text{hi}}(s_t, g)$   $\triangleleft$  get subgoal
17:        $\{wpt_{t,1}, wpt_{t,2}, \dots, wpt_{t,k}\} \leftarrow \text{Dijkstra's algorithm}(s_t, sg_t)$   $\triangleleft$  get waypoints
18:       current waypoint index  $n = 1$ 
19:     end if
20:     if achieved  $wpt_{t,n}$  or tried more than  $Dist(wpt_{t,n-1}, wpt_{t,n})$  to achieve  $wpt_{t,n}$  then
21:       current waypoint index  $+= 1$ 
22:     end if
23:      $a_t \leftarrow \pi_{\phi_1}^{\text{lo}}(s_t, wpt_{t,n+1})$   $\triangleleft$  get low-level ac-
    tion
24:   end if
25:   Env.step( $a_t$ )
26:   Train low-level agent  $Q_{\text{critic}, \theta_1}^{\text{lo}}, Q_{\text{graph}, \theta_2}^{\text{lo}}$  and  $\pi_{\phi_1}^{\text{lo}}$ , high-level agent  $Q_{\theta_3}^{\text{hi}}$  and  $\pi_{\phi_2}^{\text{hi}}$   $\triangleleft$  FLR is
    utilized when training low-level agents.
27:   if  $\tau$  % graph update freq = 0 then
28:     Update Graph  $\mathbf{G}(\mathbf{V}, \mathbf{E})$  using farthest point sampling (FPS) (Arthur & Vassilvitskii
    2007) algorithm
29:   end if
30: end for

```

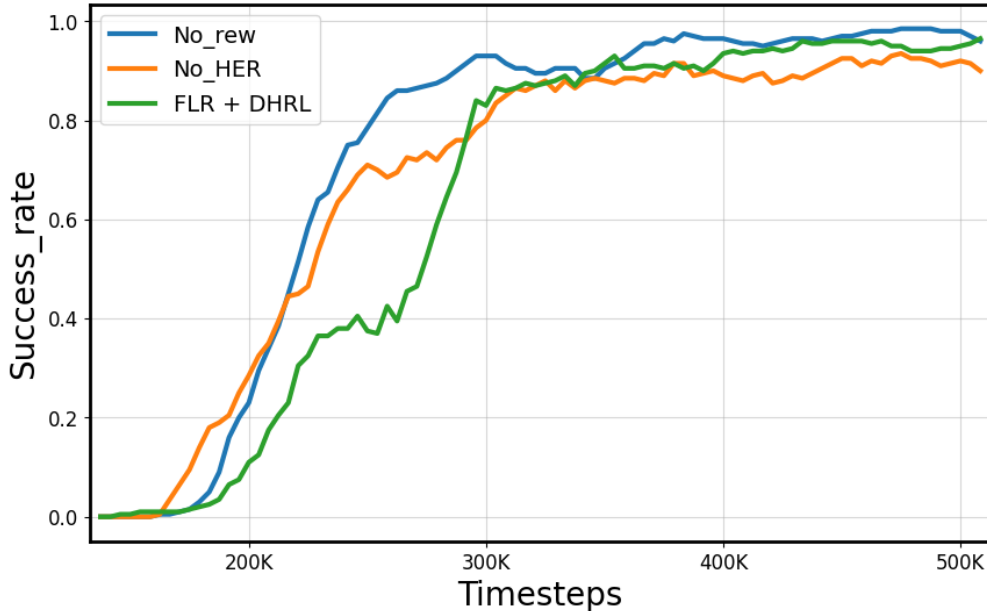


Figure 6: Comparison with oracle-based approach

B ADDITIONAL RESULTS

B.1 COMPARED WITH ORACLE

In this section, we compare our FLR+DHRL approach with an Oracle-based DHRL method in the context of the specific problem domain we address. We introduce two baseline methods for comparison: 'No rew' and 'No HER'.

- No rew: This approach assigns a reward of -1 regardless of whether the agent reaches the goal or not, specifically penalizing the occurrence of falling states. This design allows for a strong penalty for falling states.
- No HER: The 'No HER' method does not employ hindsight experience replay for falling states in the entire trajectory. This strategy ensures that the low-level policy does not learn trajectories leading to falling states, naturally encouraging the agent to avoid falling.

As shown in Figure 6, our FLR+DHRL approach does not precisely match the performance of the oracle-based method but demonstrates nearly equivalent performance. We attribute this observation to the nature of FLR, which penalizes falling states more strongly but does not explicitly incorporate mechanisms to avoid learning from undesirable transition data, akin to the 'No HER' method.

This analysis suggests that while imposing penalties for irreversible states is crucial for addressing our specific problem, avoiding such states altogether may also serve as a viable strategy for problem resolution. The FLR+DHRL approach strikes a balance between penalizing falling states and encouraging the avoidance of undesirable transitions, contributing to its effectiveness in the considered problem domain.

B.2 HYPERPARAMETERS

We utilized the majority of hyperparameters exactly as employed in the previous baselines. We experimented with varying the number of landmarks, initial rollout, and the future step in HER, as detailed in Table 12 and 3. We also report our hyperparameters in Table 4.

Table 1: Common hyperparameters setting of HGRL baselines

	PIG	DHRL	HIGL
RL algorithm	DDPG	TD3	TD3
high-level τ	-	0.005	0.005
π_h learning rate	-	0.0001	0.0001
$Q_{\pi_h}^h$ lr	-	0.001	0.001
γ_h	-	0.99	0.99
high-level train freq	-	10	10
low-level τ	0.01	0.005	0.005
π_l learning rate	0.0002	0.0001	0.0001
$Q_{\pi_l}^l$ learning rate	0.0002	0.001	0.001
γ^l	0.99	0.99	0.95
hidden layer	-	(256,256,256)	(300,300)
number of landmarks	400	300	20-150
number of novelty landmarks	-	-	20-50
batch size	128	1024	128
hindsight relabelling ratio	0.8	0.8	-
hindsight future steps	150-200	150	-

Table 2: Hyperparameters for DHRL

	DHRL
number of random rollout episodes	100-200
initial episodes without graph planning	75
gradual penalty transition rate	0.2
high-level train freq	10
target update freq	10
actor update freq	2

Table 3: Hyperparameters for PIG

	PIG
hidden layer for actors	(400,400,400,400)
hidden layer for critics	(400,400,400,400,400)

Table 4: Hyperparameters for FLR

	FLR
reward scale α	0.1
next subgoal distance ε	2
the number of sampling next subgoal	10