

1 **Response to Reviewer 1:** Thank you for the thoughtful and inspiring comments.

2 **Q1.** How much of an effect does low-level skill initialization scheme have on performance?

3 We test random skill initialization in Ant Maze task. Even random initialization is better than the non-hierarchical
4 method TRPO, shown in Fig.5 below. And more reasonable initialization results in better performance. We will explore
5 the effects of other initialization schemes in future works.

6 **Response to Reviewer 2:** Thank you for the detailed comments.

7 **Q1.** Consider potential based reward shaping in main evaluation.

8 We design a heuristic potential as the negative L2 distance between agent
9 and goal. The curve of potential reward is higher than TRPO, but significantly
10 worse than HAAR, shown in Fig.5. We note that the potential reward shaping
11 method could not take advantage of hierarchical structure; it also heavily depends on the potential function design.

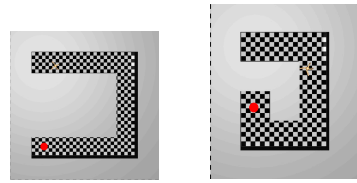


Figure 1: Task 1 Figure 2: Task 2

12 **Q2.** More robust and systematic transfer experiments.

13 We design more new tasks in Fig.1(bigger maze), 2(spiral maze) to explore the effectiveness of low-level skill transfer.
14 In the new tasks, the skill of turning right learned in Ant Maze (a) can be very useful, so low-level policy transfer shows
15 much efficiency (shown in Fig.3, 4). New tasks and the old task share much information on the high level, so “transfer
16 both” performs well, which partly comes by virtue of the transferable state representation (Konidaris, 2007).

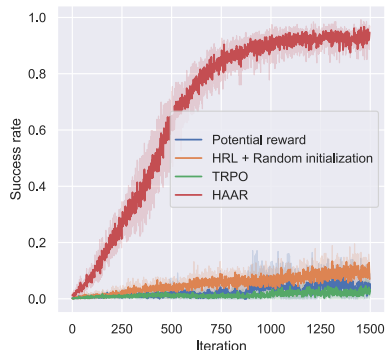
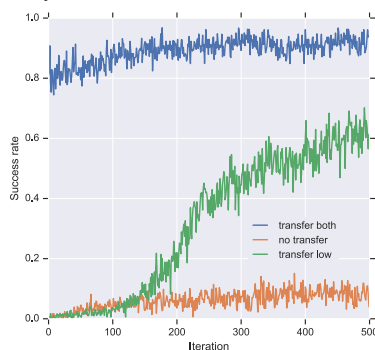


Figure 3: Transfer results for task 1.

Figure 4: Transfer results for task 2.

Figure 5: Potential based reward shaping and random low-level skill initialization in Ant Maze.

17 **Q3.** All assumptions for the proof of monotonic improvement.

18 (1) All assumptions for TRPO, including that the high-level and low-level states are Markovian; (2)The high level
19 policy is fixed while optimizing the low level policy and vice versa [lines 116-7]; (3) Discount factors $\gamma_h \rightarrow 1, \gamma_l \rightarrow 1$
20 [lines 132-3].

21 **Q4.** Experiment setting seems to be non-Markovian [line 224-5], different states may have very similar representation.

22 In experiment, the agent uses a total of 20 rays to “see” the surroundings, and the goal can always be seen regardless of
23 walls. We believe this is sufficient to distinguish between states, so it is approximately Markovian.

24 **Q5.** The benchmark, SNN4HRL, seems to run much faster in its original paper compared to in this paper.

25 The reviewer may have misread the experiment settings. Our result is actually consistent with the original SNN4HRL
26 paper. In Swimmer Maze, the numbers of samples per iteration are different in SNN4HRL paper and our paper [line
27 405]. The performance in terms of samples is consistent. For Ant Maze, SNN4HRL paper does not provide results.

28 **Q6.** A precise description of the advantage function.

29 Our definition of advantage is consistent with the conventional definition, $Q(s, a) - V(s)$. Using a one-step expansion
30 of Q , we can write it as $A_h(s^h) = E_{s_{t+k}^h \sim (\pi_h, \pi_l)} [r_t^h + \gamma_h V_h(s_{t+k}^h) | a_t^h = a^h, s_t^h = s^h] - V_h(s^h)$.

31 **Q7.** How to determine what to include in low-level state in experiments?

32 Our decision of what is included in low-level state is the same as SNN4HRL paper (such that the representation requires
33 minimal domain knowledge in the pre-training phase), described in “Problem Statement” of their paper.

34 **Response to Reviewer 3:** Thank you for the thoughtful comments.

35 **Q1.** How will the algorithm perform when starting with random low-level policies?

36 We run experiments with random initial low-level policies in Ant Maze. Results in Fig.5 show that it performs better
37 than TRPO. As expected, more meaningful low-level policies result in better performance.

38 **Q2.** Low level Markovness is not clear. Discuss the Markovness of states.

39 States for both high level and low level are Markovian. We concatenate the agent state and the high-level action a^h as
40 the low-level state, so low-level policy is still running on an MDP.