

1 We thank all reviewers for detailed and valuable comments, and will revise the paper accordingly as described below.

2 **Minor changes and typos.** We thank all reviewers for pointing those out, and will do corrections in the revision.

3 **R1 & R2: Discussions of deterministic MDPs in theoretical results.** Although our theoretical results are derived in
4 the context of deterministic MDPs, they are instructive for the practical algorithm design in general cases, and this
5 assumption is also exploited by some prior work that investigate distance metrics in MDPs (e.g. reference [11] in the
6 paper and the work of Castro mentioned by R1). Empirically, we have also showed that our method is robust to certain
7 kinds of stochasticity and outperforms the best baseline in stochastic environments (see line 40 – 48 and Figure 1).

8 **R1: Subgoals as short-term changes v.s. k -step adjacency constraint.** Interpreting the subgoals as short-term
9 changes amounts to setting a general constraint in the raw state space (e.g. HIRO) or in an embedding space (e.g.
10 FeUdal networks). However, the Euclidean distance between states in these spaces instead of the adjacency space may
11 not indicate the real adjacency relation: e.g. consider a grid-world environment where two states are separated by a wall.
12 Also, these general constraints control the maximum magnitude of the changes by a human-specified hyperparameter,
13 which is hard to choose in a principled way, while our method can learn the constraint automatically.

14 **R1: Claim of the empirical results.** We agree with the reviewer and will change the wording in the revision.

15 **R1: Comparison with two recent works.** (1) Castro proposes an algorithm for computing bisimulation metrics, which
16 reflect behavioral equivalence between states, using sampled transitions. Compared to our work, bisimulation metrics
17 depend on both the dynamics and the rewards, while the shortest transition distance depends only on the dynamics
18 and therefore can be easily applied to multi-task settings. (2) Khetarpal et al. present a theory and an algorithm of
19 affordances in RL, formulating the fact that certain states only enable certain actions. They construct the affordances
20 based on indents, i.e. desired state distributions, which are specified by humans a priori. In contrast, our method learns
21 subgoals, which can be interpreted as a kind of temporally extended indents, by RL rather than human prior.

22 **R2: Comparison with Options and EigenOptions with successor representation (SR).** (1) The Option framework
23 maintains a finite set of low-level Options (macro-actions), while our method (which falls into goal-conditioned HRL)
24 maintains an universal goal-reaching low-level policy whose behavior is modulated by subgoals. As shown in the
25 HIRO paper, goal-conditioned HRL often yields better performance than HRL with Options. (2) SR can be used to
26 measure the temporal distance between states and discover eigenoptions. Compared to our work, SR depends on both
27 environmental dynamics and a specific policy, while the shortest transition distance relies only on the dynamics.

28 **R2: Experimental settings.** We followed HIRO to use off-policy TD3 in continuous control tasks; in discrete control
29 tasks we found that whether using on-policy or off-policy methods (e.g. double DQN) does not make much difference.

30 **R3: Comparison with graph-based methods (missing baselines).** We believe that our method has essential difference
31 with graph-based methods: our method models high-level learning as a RL process and thus can tackle more general
32 problems, while current graph-based methods derive high-level policy *without* a training process, exploiting specific
33 problem structure. E.g. all graph-based works cited in the review obtain the subgoal sequence by solving a shortest-path
34 problem to a *known* goal node in a high-level graph (e.g. using Dijkstra’s algorithm), which cannot be applied to
35 more general problems where there does not exist a single “goal” state (e.g. PointGather) or the goal state needs to be
36 explored by the agent instead of being given in advance (e.g. KeyChest), as in many of our experiments. To the best of
37 our knowledge, none of these graph-based works has reported results on these general problems without additional
38 task-specific prior. Therefore, we found it hard to fairly compare our method with graph-based methods in experiments
39 and thus did not add them to baselines. In the revision, we will add these discussions to the related work section.

40 **R3: Empirical study in stochastic MDPs.** To empirically verify the stochasticity robustness
41 of our method, we have applied HRAC to a set of stochastic AntMaze tasks, which have
42 relatively larger state (30-d continuous) and action space (8-d continuous) than the Maze task.
43 We added Gaussian noise with different STDs (σ) to the (x, y) position of the ant robot at every
44 step. As shown in Figure 1, HRAC achieves similar asymptotic success rates with different
45 noise magnitudes. Due to the time limit, we only compare HRAC with the best baseline HIRO
46 when $\sigma = 0.1$, representing the noise magnitude that approximately equals to 20% of the
47 maximum step size of the ant robot on average, where HRAC achieves better performance than
48 HIRO. We will add more experimental results in stochastic environments to the revised paper.

49 **R3: Difference between HRAC and NegReward.** HRAC plugs the adjacency loss as an extra
50 term into the loss of a specific RL algorithm, rather than treat it as a negative reward.

51 **R4: Explanation of the approximation in Eq.(9).** In Eq.(9), we apply a minimizing operation instead of an
52 expectation operation over a finite policy set to approximate the original minimizing operation over the whole (indefinite)
53 policy set. Therefore, we do not require that these policies are drawn i.i.d..

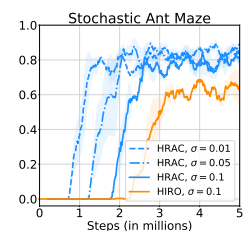


Figure 1: Learning curves.