
Emergence of a Symbolic Goal Representation with an Intelligent Tutoring System based on Intrinsic Motivation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Goal representation affects the performance of Hierarchical Reinforcement Learning (HRL) algorithms by decomposing complex problems into easier subtasks. 2 Recent studies show that representations that preserve temporally abstract environment dynamics are successful in solving difficult problems with theoretical 3 guarantees for optimality. These methods however cannot scale to tasks where 4 environment dynamics increase in complexity. On the other hand, other efforts 5 have tried to use spatial abstraction to mitigate the previous issues. Their 6 limitations include scalability to high dimensional environments and dependency on 7 prior knowledge. In this work, we propose a novel three-layer HRL algorithm that 8 introduces, at different levels of the hierarchy, both a spatial and a temporal goal 9 abstraction. We provide a theoretical study of the regret bounds of the learned 10 policies. We evaluate the approach on complex continuous control tasks, demonstrating 11 the effectiveness of spatial and temporal abstractions learned by this approach. 12 13

14 1 Introduction

15 Open-ended learning, as defined by Doncieux et al. [2018], should be solved by an unsupervised 16 acquisition of a hierarchy of adapted representations. Along with hierarchical descriptions of actions 17 in neuroscience (eg Grafton and de C. Hamilton [2007]) and in behavioural psychology (eg. Eckstein 18 and Collins [2021]), in machine learning, *Hierarchical Reinforcement Learning* (HRL) [Barto and 19 Mahadevan, 2003] tackles the environment’s complexity by introducing a hierarchical structure of 20 agents that work at different levels of temporal and behavioral abstractions.

21 Recent works (Vezhnevets et al. [2017], Nachum et al. [2019], Zhang et al. [2023], Li et al. [2021]) 22 have shown that an abstract goal representations can propose semantically meaningful subgoals and to 23 solving more complex tasks. In particular, representations that capture environment dynamics over an 24 abstract temporal scale have been shown to provide interesting properties with regards to bounding the 25 suboptimality of learned policies under abstract goal spaces (Nachum et al. [2019], Abel et al. [2020]), 26 as well as efficiently handling continuous control problems. However, temporal abstractions that 27 capture aspects of the environment dynamics (Ghosh et al. [2019], Savinov et al. [2019], Eysenbach 28 et al. [2019], Zhang et al. [2023], Li et al. [2021]) still cannot scale to environments where the 29 pairwise state reachability relation is complex. This situation typically occurs when temporally 30 abstract relations take into account more variables in the state space. The main limitation of these 31 approaches is the lack of a spatial abstraction to generalise such relations over states.

32 Alternatively, other works (Kulkarni et al. [2016], Illanes et al. [2020], Garnelo and Shanahan [2019] 33) have studied various forms of spatial abstractions for goal spaces. These abstractions effectively 34 group states with similar roles in sets to construct a discrete goal space. The advantage of such 35 representation is a smaller size exploration space that expresses large and long-horizon tasks. In 36 contrast to these algorithms that require varying levels of prior knowledge, GARA Zadem et al.

37 [2023] gradually learns such spatial abstractions by considering reachability relations between sets
38 of states. We refer to this abstraction as *reachability-aware abstraction*. While such representation
39 is efficient for low-dimensional tasks, scalability remains an issue due to the lack of a temporal
40 abstraction mechanism. What is challenging about scaling GARA’s approach to more complex
41 environments is exactly what makes the set-based representation effective: the low-level agent must
42 learn how to reach a set of states that, especially in the initial phases of the algorithm when the
43 abstraction is coarser, may be “far” away. We tackle this problem introducing a new agent in the
44 hierarchy that introduces a *temporal abstraction*. Such an agent learns to select intermediate subgoals
45 that: can be reached from a state s executing the low-level agent; and helps constructing a trajectory
46 from s to a goal abstract state.

47 Here, we propose a three-layer HRL algorithm achieving both a temporal and spatial abstraction to
48 capture the environment dynamics. We motivate the use of temporal abstraction as the key factor
49 that can scale the abstraction proposed in Zadem et al. [2023], and the reachability-aware spatial
50 abstraction as a way to efficiently represent goals in complex tasks. We complement the approach by
51 adding theoretical guarantees on the bounds of suboptimality of policies learned under this abstraction.
52 Our approach is empirically evaluated on a set of challenging continuous control tasks.

53 2 Related Work

54 Building on ideas for introducing hierarchy in Reinforcement Learning [Sutton et al., 1999, Barto and
55 Mahadevan, 2003, Dayan and Hinton, 1992], recent advances have managed to considerably elevate
56 HRL algorithms to tackle complex continuous control environments. For instance, Nachum et al.
57 [2018] introduces a two-level hierarchy that sample goals from a pre-defined oracle on the state space.
58 This approach provides a good basis for HRL algorithms as it is generic and addresses non-stationary
59 learning but may still be suboptimal as the goal sampling is unconstrained in the oracle. To remedy
60 this, Ghosh et al. [2019], Savinov et al. [2019], Eysenbach et al. [2019], Zhang et al. [2023], Li et al.
61 [2021] learn different goal representations that try to capture the environment dynamics. This idea
62 has been validated under different theoretical formulations [Nachum et al., 2019, Abel et al., 2020, Li
63 et al., 2021]. In particular, Li et al. [2021] learns a latent representations based on slow-dynamics in
64 the state space. The idea is that meaningful temporally abstract relations (over k steps) are expressed
65 by state features that slowly change over time. However, these features may not be always sufficient
66 to capture all the critical information about dynamics. Both Savinov et al. [2019] and Zhang et al.
67 [2023] use k -step reachability relations as a characterisation for environment dynamics. Their idea is
68 to learn if goals (mappings of states in an embedding / oracle) reach a potential goal in k steps. These
69 relations are later used to drive the sampling of goals that can be reached, resulting in more efficient
70 learning. However, such learned pairwise relations are binary and lack the information of which
71 goals can be reached by applying a specific policy. Additionally, without any spatial abstraction, it
72 can be difficult to learn these relations for a complex transition relation (e.g. a relation that require
73 monitoring more that few variables).

74 To introduce spatial abstraction, we study the work by Zadem et al. [2023] in which the authors
75 introduce GARA, a spatial abstraction for the goal space that captures richer information from k
76 step reachability relations. This algorithm progressively learns a discretisation of the state space that
77 serves as an abstract goal space. This abstraction generalizes reachability relations over sets of states,
78 greatly reducing the difficulty of the learning problem. This approach however was only validated on
79 low dimensional environments and suffers from scalability issues. As GARA starts learning from
80 a coarse abstraction (composed of a small number of large sets), the goal set is often distant from
81 the current state, thus it can be difficult to learn meaningful policies that manage to reach a desired
82 goal set. Under such circumstances, the abstraction cannot be refined as it lacks any mechanism to
83 propose easier, more granular subgoals. To alleviate this discrepancy, we introduce a new agent in
84 the hierarchy of GARA that applies a *temporal abstraction* [Sutton et al., 1999].

85 3 Theoretical properties of the refinement

86 We motivate the adoption of the goal-space abstraction and the reachability-aware refinement showing
87 that: (i) there exists a bound on the suboptimality of policies trained with a reachability-aware
88 abstraction; and (ii) the reachability-aware refinement gradually finds a reachability-aware abstraction.
89 The theoretical results hold under the assumptions that the environment M is *deterministic* and the
90 reward signal r_{ext} is bounded in the environment.

91 **Suboptimality bounds:** Since the introduction of a set-based abstraction changes how the hierar-
 92 chical policies are learned, we are interested in measuring how much the optimal policy $\pi_{\mathcal{N}}^*$ (learned
 93 on the reachability-aware abstraction \mathcal{N}) deviates from the optimal hierarchical policy π^* (learned
 94 without abstraction). Each policy’s performance is characterised by a value function measuring its
 95 future cumulative rewards. In other words, the comparison of both policies comes down to bounding
 96 the difference between the value functions. This difference will be computed on the states visited by
 97 each optimal policy. These states are referred to as the trajectory of a policy.

98 Our results show that there exists a bound on how much the value function under $\pi_{\mathcal{N}}^*$ deviates from
 99 the value function under π^* . We also observe that if the size of the abstract goals is unbounded then
 100 the deviation between the two policies can progressively with time t , reaching a peak during the
 101 middle point of the trajectory, before narrowing and eventually converging near its end. On the other
 102 hand, if we assume that the size of goal sets is bounded, then a tighter and more stable bound is
 103 obtained and the deviation between policies does not depend on the timestep t . Please refer to Annex
 104 C for the details and theorems.

105 **Reachability-aware abstraction through iterative refinement:** To complement the characteri-
 106 sation of an upper bound on the suboptimality of $\pi_{\mathcal{N}}^*$, we show that such an abstraction is indeed
 107 obtainable in a finite number of iterative refinements as proposed in Zadem et al. [2023]. Please refer
 108 to Annex C for the theorem.

109 4 Spatio-Temporal Abstraction via Reachability

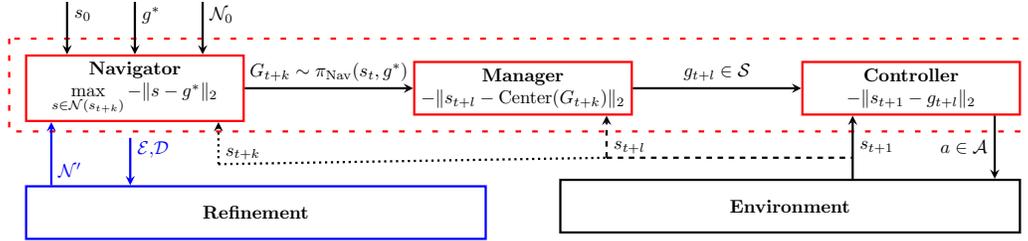


Figure 1: **Architecture of STAR.** The algorithm’s inputs are the initial state s_0 , the task goal g^* , and an initial abstraction \mathcal{N}_0 . STAR runs in a feedback loop a Feudal HRL algorithm (dashed red block) and an abstraction refinement (blue box). The solid red blocks show the HRL agents (*Navigator*, *Manager*, *Controller*) and the agents reward in the bottom. The agents run at different timescales ($k > l > 1$), shown with the solid, dashed, and dotted lines carrying the feedback from the environment to the agents. The Refinement uses as inputs the past episodes (\mathcal{D}) and the list of abstract goals (\mathcal{E}) visited during the last episode, and outputs an abstraction. g^* is the task’s goal.

110 We propose a HRL algorithm, Spatio-Temporal Abstraction via Reachability (STAR), that learns, at
 111 the same time, a spatial goal abstraction \mathcal{N} and policies at multiple time scales (See Algorithm 1 in
 112 the Annex). The STAR algorithm, shown in Figure 1, has two main components: a 3-levels Feudal
 113 HRL algorithm (enclosed in the red dashed lines); and an abstraction refinement component (shown
 114 in the blue solid lines). STAR runs the Feudal HRL algorithm and the abstraction refinement in a
 115 feedback loop, refining the abstraction \mathcal{N} at the end of every learning episode. The feudal architecture
 116 of STAR is composed of a hierarchy with three agents:

- 117 1. *Navigator*: the highest-level agent samples, every k steps, an abstract **goal** $G \in \mathcal{G}$ that
 118 should help to reach the task goal g^* from the current agent’s state ($G_{t+k} \sim \pi_{\text{Nav}}(s_t, g^*)$).
- 119 2. *Manager*: the mid-level agent is conditioned by the navigator goal G and every k steps,
 120 picks **subgoals** in the state space ($g_{t+l} \sim \pi_{\text{Man}}(s_t, G_{t+k})$). The first intuition is that, for a
 121 possibly very far goal G , the *Manager* samples an intermediate subgoal $g \in \mathcal{S}$ of a difficulty
 122 level adapted to the current non-optimal policy. The *Manager* receives an intrinsic reward
 123 expressing the distance between the *Controller*’s state s_t and the target abstract state G_{t+k} :

$$r_{\text{Manager}}(s_t, G_{t+k}) := -\|s_t - \text{Center}(G_{t+k})\|_2,$$

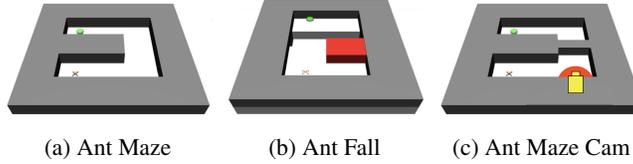


Figure 2: Ant environments

124 where $\text{Center}(G_{t+k})$ is the center of the goal G_{t+k} . Specifically, a set is a hyperrectangle
 125 (i.e., a multi-dimensional interval), so the center’s coordinates are the mean between the min
 126 and max of the interval.

127 3. *Controller*: the low-level policy is goal-conditioned by the *Manager*’s subgoal g and samples
 128 actions to reach given goal ($a \sim \pi_{\text{Cont}}(s_t, g_{t+l})$).¹

129 The algorithm learns the abstraction iteratively. Every refinement obtains a finer abstraction \mathcal{N}'
 130 from \mathcal{N} . Intuitively, \mathcal{N}' will split at least a goal $G_1 \in \mathcal{G}_{\mathcal{N}}$ in two goals $G'_1, G''_1 \in \mathcal{G}_{\mathcal{N}'}$ if there are
 131 different states in G_1 (i.e., $s_a, s_b \in G_1$) that cannot reach the same target $G_2 \in \mathcal{G}$ when applying the
 132 same low-level policy. We define such *reachability property* precisely in Annex B, intuitively the
 133 refinement separates goal states that "behave differently" under the same low level policy (i.e., \mathcal{N}'
 134 would represent more faithfully the environment dynamics). See details in Annex B.

135 5 Experimental Evaluation

136 We answer the following research questions: 1) Do the spatial and temporal abstraction of STAR
 137 allow for more data-efficient learning? 2) Does the reachability-aware abstraction scale to more
 138 complex environments compared to a more concrete reachability relation?

139 5.1 Environment Setup

140 We evaluate our approach on a set of challenging tasks in the Ant environments (Fig.2) adapted from
 141 Duan et al. [2016] and popularised by Nachum et al. [2018]. We propose the following tasks:

142 1. **Ant Maze**: in this task, the ant must navigate a 'D'-shaped maze to reach the exit positioned at the
 143 top left. 2. **Ant Fall**: the environment is composed of two raised platforms separated by a chasm.
 144 The ant starts on one of the platforms and must safely cross to the exit without falling. A movable
 145 block can be pushed into the chasm to serve as a bridge. Besides the precise maneuvers required by
 146 the ant, falling into the chasm is a very likely yet irreversible mistake. 3. **Ant Maze Cam**: this is a
 147 more challenging version of Ant Maze. The upper half of the maze is fully blocked by an additional
 148 obstacle that can only be opened when the ant looks at the camera (in yellow in Fig. 2c) when on the
 149 red spot. The exit remains unchanged.

150 5.2 Comparative Analysis

151 We compare STAR with the following algorithms: 1. **GARA** [Zadem et al., 2023] 2. **HIRO**
 152 [Nachum et al., 2018] 3. **HRAC** [Zhang et al., 2023] 4. **LESSON** [Li et al., 2021]

153 In line with HIRO and HRAC, STAR relies on an oracle $\psi(s)$ that transforms the observations of
 154 the high-level agents (Manager and Navigator). In practice $\psi()$ corresponds to a feature selection
 155 applied to states. In contrast, LESSON learns a latent goal space without an oracle. In Ant Maze
 156 $\psi(s) = (x, y)$, in Ant Fall $\psi(s) = (x, y, z)$, and in Ant Maze Cam, $\psi(s) = (x, y, \theta_x, \theta_y, \theta_z)$.

157 Fig.4 shows that STAR outperforms all of the state-of-art approaches by reaching a higher success
 158 rate with less timesteps. In particular GARA, operating only under a spatial abstraction mechanism
 159 is unable to solve Ant Maze, the easiest task in this analysis. HIRO on the other hand learns less
 160 efficient policies due to it lacking a spatial abstraction component. These results show that STAR,
 161 which combines temporal and spatial abstractions, is a more efficient approach.

162 To discuss the second research question, we first observe that, while the high-level dynamics of
 163 Ant Maze can be captured by the x, y dimensions, the dynamics of Ant Fall require all the x, y, z
 164 dimensions (z expresses if the ant is safely crossing above the pit or if it has fallen), and Ant Maze

¹In the following, we use the upper-case G letter for goals in \mathcal{G} and the lower-case g for subgoals in \mathcal{S} .

165 Cam requires x, y, θ_x, θ_y , and θ_z (the orientation angles are necessary to unlock the access to the
 166 upper part of the maze). Fig.4 shows that HRAC is unable to capture meaningful relations between
 167 subgoals and fails at solving either Ant Fall or Ant Maze Cam due to the increased complexity
 168 in capturing the high-level task dynamic. Similarly, LESSON is unable to learn a good subgoal
 169 representation in Ant Maze Cam using slowness principle dynamics since it doesn't apply to $\theta_x, \theta_y, \theta_z$.
 170 Instead, STAR is capable of abstracting these dimensions and converging to a successful policy.

171 **5.3 Representation Analysis**

172 We answer the third research question examining the progress of the STAR's navigator at different
 173 timesteps during learning when solving the Ant Maze and Ant Fall tasks (visualizing the 5 dimensional
 174 representation for Ant Maze Cam is less feasible). In Fig.3, we plot the frequency of the goals visited
 175 by the navigator when evaluating a policy learned after 1M, 2M, and 3M timesteps (we average the
 176 frequency over 5 different evaluations of 500 maximum timesteps each).

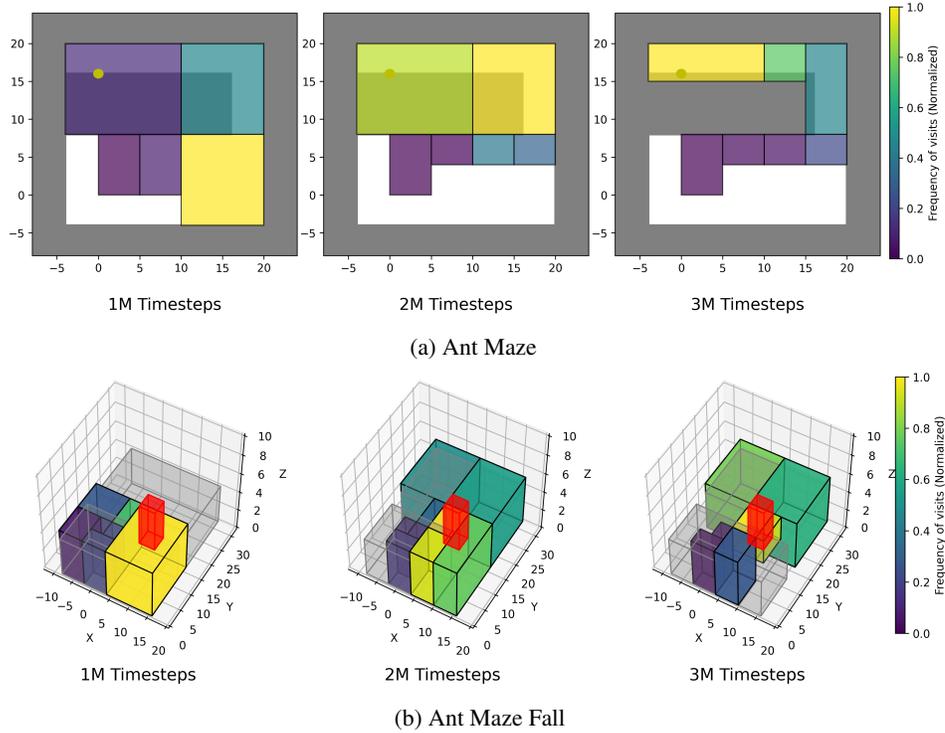


Figure 3: Frequency of goals visited by the navigator when evaluating a policy learned after 1M, 2M, and 3M timesteps (averaged over 5 different evaluations with 500 maximum timesteps). The subdivision (squares for Ant Maze, boxes for Ant Maze Fall) represent (abstract) goals. The color represents the frequency of visits of each goal. Grey areas correspond to the obstacles of the environment in Ant Maze,

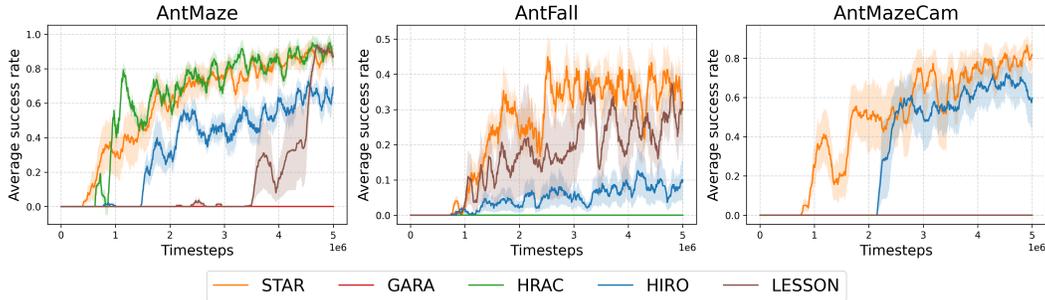


Figure 4: Comparative evaluation on Ant environments (Averaged over 5 runs)

References

- 177
178 David Abel, Nate Umbanhowar, Khimya Khetarpal, Dilip Arumugam, Doina Precup, and Michael Littman.
179 Value preserving state-action abstractions. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of*
180 *the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings*
181 *of Machine Learning Research*, pages 1639–1650. PMLR, 26–28 Aug 2020. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v108/abel20a.html)
182 [press/v108/abel20a.html](https://proceedings.mlr.press/v108/abel20a.html).
- 183 Andrew G. Barto and Sridhar Mahadevan. Recent advances in hierarchical reinforcement learning. *Discrete*
184 *Event Dynamic Systems*, 13(1), 2003.
- 185 Peter Dayan and Geoffrey E Hinton. Feudal reinforcement learning. In *NeurIPS*, volume 5, 1992.
- 186 Stephane Doncieux, David Filliat, Natalia Díaz-Rodríguez, Timothy Hospedales, Richard Duro, Alexandre
187 Coninx, Diederik M. Roijers, Benoît Girard, Nicolas Perrin, and Olivier Sigaud. Open-ended learning: A
188 conceptual framework based on representational redescription. *Frontiers in Neurobotics*, 12, sep 2018.
- 189 Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement
190 learning for continuous control. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of*
191 *The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning*
192 *Research*, pages 1329–1338, New York, New York, USA, 20–22 Jun 2016. PMLR. URL [https://proceedings.](https://proceedings.mlr.press/v48/duan16.html)
193 [mlr.press/v48/duan16.html](https://proceedings.mlr.press/v48/duan16.html).
- 194 Maria K Eckstein and Anne G E Collins. How the mind creates structure: Hierarchical learning of action
195 sequences. In Cognitive Science Society, editor, *CogSci Conference of the Cognitive Science Society*,
196 volume 43, pages 618–624, 2021.
- 197 Ben Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. Search on the replay buffer: Bridging planning and
198 reinforcement learning. In *NeurIPS 2019*, 2019.
- 199 Marta Garnelo and Murray Shanahan. Reconciling deep learning with symbolic artificial intelligence: represent-
200 ing objects and relations. *Current Opinion in Behavioral Sciences*, 29:17–23, 2019. ISSN 2352-1546. doi:
201 <https://doi.org/10.1016/j.cobeha.2018.12.010>. Artificial Intelligence.
- 202 Dibya Ghosh, Abhishek Gupta, and Sergey Levine. Learning actionable representations with goal conditioned
203 policies. In *ICLR 2019*, 2019.
- 204 Scott T. Grafton and Antonia F. de C. Hamilton. Evidence for a distributed hierarchy of action representation in
205 the brain. *Human Movement Science*, 26(4):590–616, 2007.
- 206 León Illanes, Xi Yan, Rodrigo Toro Icarte, and Sheila A. McIlraith. Symbolic plans as high-level instructions for
207 reinforcement learning. In *AAAI*, 2020.
- 208 Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical deep reinforcement
209 learning: Integrating temporal abstraction and intrinsic motivation. In *NeurIPS*, volume 29, 2016.
- 210 Siyuan Li, Lulu Zheng, Jianhao Wang, and Chongjie Zhang. Learning subgoal representations with slow
211 dynamics. In *ICLR*, 2021.
- 212 Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning.
213 In *NeurIPS 2018*, 2018.
- 214 Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Near-optimal representation learning for
215 hierarchical reinforcement learning. In *ICLR*, 2019.
- 216 Nikolay Savinov, Anton Raichuk, Damien Vincent, Raphaël Marinier, Marc Pollefeys, Timothy P. Lillicrap, and
217 Sylvain Gelly. Episodic curiosity through reachability. In *ICLR*, 2019.
- 218 Richard S. Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal
219 abstraction in reinforcement learning. *Artificial Intelligence*, 112, 1999.
- 220 Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and
221 Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. *CoRR*, abs/1703.01161, 2017.
- 222 Mehdi Zadem, Sergio Mover, and Sao Mai Nguyen. Goal space abstraction in hierarchical reinforcement
223 learning via set-based reachability analysis, 2023.
- 224 Tianren Zhang, Shangqi Guo, Tian Tan, Xiaolin Hu, and Feng Chen. Adjacency constraint for efficient
225 hierarchical reinforcement learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):
226 4152–4166, 2023. doi: 10.1109/TPAMI.2022.3192418.

Algorithm 1 STAR

Input: Learning environment E .
Output: Computes π_{Nav} , π_{Man} and π_{Cont}

- 1: $\mathcal{D}_{\text{Navigator}} \leftarrow \emptyset, \mathcal{D}_{\text{Manager}} \leftarrow \emptyset, \mathcal{D}_{\text{Controller}} \leftarrow \emptyset, \mathcal{G} \leftarrow \{\mathcal{S}\}$
- 2: **for** $t \leq \text{max_timesteps}$ **do**
- 3: $\mathcal{E} \leftarrow \emptyset$
- 4: $s_{\text{init}} \leftarrow$ initial state from $E, s_t \leftarrow s_{\text{init}}$
- 5: $G_s \leftarrow G \in \mathcal{G}$ such that $s_t \in G$
- 6: $G_d \sim \pi_{\text{Nav}}(G_s, g^*)$
- 7: $g_t \sim \pi_{\text{Man}}(s_t, G_d)$
- 8: **while** true **do**
- 9: $\mathcal{E} \leftarrow \mathcal{E} \cup \{(G_s, G_d)\}$
- 10: $a_t \sim \pi_{\text{Cont}}(s_t, g_t)$
- 11: $(s_{t+1}, r_t^{\text{ext}}, \text{done}) \leftarrow$ execute the action a_t at s_t in E
- 12: $r_{\text{Controller}} = -\|g_t - s_t\|_2$
- 13: Update π_{Cont}
- 14: **if** not done **then**
- 15: $s_t \leftarrow s_{t+1}, t \leftarrow t + 1$
- 16: **if** $t \bmod l = 0$ **then**
- 17: $\mathcal{D}_{\text{Manager}} \leftarrow (s_{t-l}, G_d, g_{t-l}, s_t, r_{\text{Manager}}, \text{done})$
- 18: Update π_{Man}
- 19: $g_t \sim \pi_{\text{Man}}(s_t, G_d)$
- 20: **if** $t \bmod k = 0$ **then**
- 21: Update $\mathcal{D}_{\text{Navigator}} \leftarrow (s_{t-k}, G_d, s_t, r_{\text{Navigator}}, \text{done})$
- 22: Update π_{Nav}
- 23: $G_s \leftarrow G \in \mathcal{G}$ such that $s_t \in G$
- 24: $G_d \sim \pi_{\text{Nav}}(s_t, g_{\text{exit}})$
- 25: **else**
- 26: Update \mathcal{F}_k with the data from $\mathcal{D}_{\text{Navigator}}$
- 27: $\mathcal{G} \leftarrow \text{Refine}(\mathcal{G}, \mathcal{E}, \mathcal{F}_k)$
- 28: break the **while** loop and start a new episode

228 **B Refining \mathcal{N} via Reachability Analysis**

229 While we follow the high level description of the refinement procedure of the GARA [Zadem et al., 2023]
 230 algorithm, we adapt it to our notation and to the new theoretical results on the refinement we present later (which
 231 holds for both GARA and STAR). Furthermore, in the rest of this Section and in Section C, we will assume a
 232 2-level hierarchy, where π_{Nav} is the high-level policy (e.g., π_{high}), and π_{low} is the hierarchical policy obtained
 233 composing π_{Man} and π_{Cont} .

234 We first define the k -step reachability relation for a goal-conditioned policy $\pi_{\mathcal{N}_{\text{Low}}}$:

$$R_{\pi_{\mathcal{N}_{\text{Low}}}}^k(G_i, G_j) := \left\{ s' \in \mathcal{S} \mid s \in G_i, s \xrightarrow[\pi_{\mathcal{N}_{\text{Low}}}(\cdot, G_j)]{k} s' \right\},$$

235 where $s \xrightarrow[\pi_{\mathcal{N}_{\text{Low}}}(\cdot, G_j)]{k} s'$ means that s can reach s' by executing the policy $\pi_{\mathcal{N}_{\text{Low}}}(\cdot, G_j)$ (targeting G_j) in k steps.

236 In other words, $R_{\pi_{\mathcal{N}_{\text{Low}}}}^k(G_i, G_j)$ is the set of states reached when starting in any state in G_i and applying the
 237 policy $\pi_{\mathcal{N}_{\text{Low}}}(\cdot, G_j)$ for k steps.

238 The algorithm uses the notion of reachability property among a pair of abstract goals:

239 **Definition 1 (Pairwise Reachability Property)** Let $\mathcal{N} : \mathcal{S} \rightarrow 2^{\mathcal{S}}$ be a set-based abstraction and $G_i, G_j \in$
 240 $\mathcal{G}_{\mathcal{N}}$. \mathcal{N} satisfies the pairwise reachability property for (G_i, G_j) if $R_{\pi_{\mathcal{N}_{\text{Low}}}}^k(G_i, G_j) \subseteq G_j$.

241 The algorithm decides to refine the abstract representation after an episode of the HRL algorithm. Let $\mathcal{E} :=$
 242 $\{G_0, \dots, G_n\}$ be the list of goals visited in the last episode. The refinement algorithm analyzes all the pairs
 243 $(G_i, G_{i+1}) \in \mathcal{E}$, for $0 \leq i < n$, and refines \mathcal{N} in a new abstraction \mathcal{N}' "splitting" G_i if it does not satisfy
 244 the pairwise reachability property. Each refinement obtains a new, finer abstraction \mathcal{N}' where the reachability
 245 property is respected in one more goal. We formalize when an abstraction \mathcal{N}' refines an abstraction \mathcal{N} with
 246 respect to the reachability property as follows:

247 **Definition 2 (Pairwise Reachability-Aware Refinement)** Let $\mathcal{N} : \mathcal{S} \rightarrow 2^{\mathcal{S}}$ and $\mathcal{N}' : \mathcal{S} \rightarrow 2^{\mathcal{S}}$ be two set-
 248 based abstractions such that there exists $G_i \in \mathcal{G}_{\mathcal{N}}$, $\mathcal{G}_{\mathcal{N}'} = (\mathcal{G}_{\mathcal{N}} \setminus \{G_i\}) \cup \{G'_1, G'_2\}$, $G'_1 \cup G'_2 = G_i$, and
 249 $G'_1 \cap G'_2 = \emptyset$. \mathcal{N}' refines \mathcal{N} (written as $\mathcal{N}' \prec \mathcal{N}$) if, for some $G_j \in \mathcal{G}_{\mathcal{N}}$, \mathcal{N}' satisfies the pairwise reachability
 250 property for (G'_1, G_j) , while \mathcal{N} does not satisfy the pairwise reachability property for (G_i, G_j) .

251 **Approximating the Reachability Property.** STAR checks the reachability property for (G_i, G_j) by
 252 approximating a reachability relation among abstract states with a forward model (a neural network) $\mathcal{F}_k : \mathcal{S} \times \mathcal{G} \rightarrow \mathcal{S}$
 253 such that $\mathcal{F}_k(s_t, G_j) \simeq s_{t+k}$. This model is trained from exploration data to approximate the
 254 reached state after applying the policy $\pi_{\mathcal{N}_{\text{Low}}}(s, G_j)$ for k steps. Checking the reachability property for (G_i, G_j)
 255 amounts to checking the output of $\mathcal{F}_k(s \in G_i, G_j)$ for all states $s \in G_i$ with a neural network reachability
 256 analysis tool [?]. If the reachability property is not valid, G_i is progressively split into smaller regions to identify
 257 subsets where the reachability property holds. Further detail can be found in Zadem et al. [2023].

258 C Theoretical Properties of the Refinement

259 In this section, we motivate the adoption of the goal-space abstraction and the reachability-aware refinement
 260 showing that: (i) there exists a bound on the sub-optimality of policies trained with a reachability-aware
 261 abstraction; and (ii) the reachability-aware refinement gradually finds a reachability-aware abstraction. Our
 262 results apply to both STAR and GARA [Zadem et al., 2023], and all the proofs are available in the Appendix ??.

263 The theoretical results hold under the assumptions that the environment M is *deterministic* and the reward
 264 signal r_{ext} is bounded in the environment. Consequently, we assume that the distance separating a state
 265 $s \in \mathcal{S}$ from all the states $s' \in \mathcal{S}$ that s can reach in one step is bounded. Thus, there is an upper bound
 266 $R_{\text{max}} := \max_{s, s' \in \mathcal{S}, \sum_{a \in \mathcal{A}} P(s' | s, a) \geq 0} \|s - s'\|_2$ on the reward signal.

267 Let π^* be the optimal hierarchical policy composed by a high-level policy $g \sim \pi_{\text{high}}^*(s, g^*)$ that samples $g \in \mathcal{S}$,
 268 and a low-level policy $a \sim \pi_{\text{low}}^*(s, g)$ that samples actions $a \in \mathcal{A}$. Since the environment is deterministic,
 269 there exists an *optimal high-level trajectory* containing the goals sampled with π_{high}^* and an *optimal low-level*
 270 *trajectory* containing all the visited states:

$$\mathcal{T}_{\text{high}}^* := \{g_0, g_1, \dots, g_m\}, \quad \mathcal{T}_{\text{Low}}^* := \{s_0, s_1, \dots, s_{m \cdot k}\}, \quad \text{with } s_{i \cdot k} = g_i, \quad \text{for } 0 \leq i \leq m.$$

271 Let $\mathcal{N} : \mathcal{S} \rightarrow 2^{\mathcal{S}}$ be a set-based abstraction. We write $\pi_{\mathcal{N}}^*$ for the *optimal* hierarchical policy obtained with
 272 the abstraction \mathcal{N} . We write $\mathcal{T}_{\mathcal{N}_{\text{High}}}^*$ and $\mathcal{T}_{\mathcal{N}_{\text{Low}}}^*$ for the optimal high- and low-level trajectories respectively.
 273 Below, we provide an upper bound on the difference between the optimal hierarchical policy π^* and the optimal
 274 hierarchical policy $\pi_{\mathcal{N}}^*$ when \mathcal{N} is a *reachability-aware*.

275 **Definition 3 (Reachability-Aware Abstraction)** Let $\mathcal{N} : \mathcal{S} \rightarrow 2^{\mathcal{S}}$ be a set-based abstraction, $\pi_{\mathcal{N}}^*$ the corre-
 276 sponding optimal hierarchical policy, and $\mathcal{T}_{\text{high}}^*$ the optimal high-level trajectory from π_{high}^* . \mathcal{N} is a reachability-
 277 aware abstraction with respect to $\mathcal{T}_{\text{high}}^*$ if:

278 1. States are contained in their abstraction: $\forall s \in \mathcal{S}, s \in \mathcal{N}(s)$.

279 2. The abstractions of the goals in the optimal trajectory are disjoint:

$$\forall g_i, g_j \in \mathcal{T}_{\text{high}}^*, (g_i \neq g_j \rightarrow \mathcal{N}(g_i) \cap \mathcal{N}(g_j) = \emptyset).$$

280 3. The abstractions of each consecutive goals in the optimal trajectory satisfy the pairwise reachability
 281 property:

$$\forall g_i, g_{i+1} \in \mathcal{T}_{\text{high}}^*, R_{\pi_{\mathcal{N}_{\text{Low}}}^*}^k(\mathcal{N}(g_i), \mathcal{N}(g_{i+1})) \subseteq \mathcal{N}(g_{i+1}).$$

282 4. The reward in the final abstract goal $\mathcal{N}(g_m)$ is bounded:

$$\exists \epsilon > 0, \forall s \in \mathcal{N}(g_m), |r_{\text{ext}}(g_m) - r_{\text{ext}}(s)| \leq \epsilon.$$

283 **Theorem 1 (Sub-optimal Learning)** Let M be a deterministic environment with task goal $g^* \in \mathcal{S}$ and
 284 $r_{\text{ext}}(s) = -\|g^* - s\|_2$. Let $\mathcal{N} : \mathcal{S} \rightarrow 2^{\mathcal{S}}$ be a reachability-aware abstraction with respect to $\mathcal{T}_{\text{high}}^*$. Then, for
 285 $s_0 \in \mathcal{T}_{\text{Low}}^*$ and $s'_0 \in \mathcal{T}_{\mathcal{N}_{\text{Low}}}^*$ we have that:

$$|V_{\pi^*}(s_0) - V_{\pi_{\mathcal{N}}^*}(s'_0)| \leq \left(\sum_{i=0}^{\frac{mk}{2}} \gamma^i i + \sum_{i=\frac{mk}{2}}^{mk} \gamma^i (mk - i) \right) \cdot 2R_{\text{max}} + \frac{1 - \gamma^{mk+1}}{1 - \gamma} \epsilon, \quad (1)$$

286 where $V_{\pi}(s)$ is the value function for a policy π [?].

287 Moreover, if there exists a $B \geq \epsilon > 0$ such that for all $1 \leq i \leq m$, $\max_{x, y \in \mathcal{N}(g_i)} \|x - y\| \leq B$, then $\forall s_i \in \mathcal{T}_{\text{Low}}^*$

288 and $\forall s'_i \in \mathcal{T}_{\mathcal{N}_{\text{Low}}}^*$ we have that:

$$|V_{\pi_{\text{Low}}^*}(s_0) - V_{\pi_{\mathcal{N}_{\text{Low}}}^*}(s'_0)| \leq \frac{1 - \gamma^{mk+1}}{1 - \gamma} (kR_{\text{max}} + B). \quad (2)$$

289 Equation (1) in the above theorem provides a bound on the sub-optimality of the hierarchical policy when
 290 trained under a set-based reachability-aware abstraction \mathcal{N} . Intuitively, the worst trajectory $\mathcal{T}_{\mathcal{N}_{\text{Low}}}^*$, starting from
 291 $s'_0 \in \mathcal{N}(s_0)$, can progressively deviate from $\mathcal{T}_{\text{Low}}^*$ as i increases. When $i \geq \frac{mk}{2}$, the trajectories progressively
 292 converge around $\mathcal{N}(g_m)$. Equation (2) defines a tighter upper bound when there is a bound B on the maximum
 293 distance between two states in each abstract goal in $\mathcal{T}_{\mathcal{N}_{\text{High}}}^*$. In practice, the existence of the bound B is valid
 294 when the state space is bounded. In this case, the deviation of the two trajectories is independent from i and is
 295 stable across time.

296 **Lemma 1** *Let \mathcal{N} and \mathcal{N}' be two set-based abstractions such that $\mathcal{N}' \prec \mathcal{N}$ and \mathcal{N} satisfies the Conditions (1),
 297 (2), and (3) (but not (4)) of a reachability-aware abstraction (Definition 3). Also, let $G_i \in \mathcal{T}_{\mathcal{N}_{\text{High}}}^*$ (note that
 298 $G_i \in \mathcal{G}_{\mathcal{N}}$), and G_i be the goal refined in \mathcal{N}' . Then, the abstraction \mathcal{N}' satisfies the following:*

- 299 1. $\exists g_i \in \mathcal{T}_{\text{High}}^*$ such that \mathcal{N} does not satisfy the reachability property for $(\mathcal{N}(g_i), \mathcal{N}(g_{i+1}))$, while \mathcal{N}'
 300 does for $(\mathcal{N}'(g_i), \mathcal{N}(g_{i+1}))$.
- 301 2. If there exists $g_j \in \mathcal{T}_{\text{High}}^*$ such that $g_j \in \mathcal{N}(g_i)$, then $g_j \notin \mathcal{N}'(g_i)$.

302 **Theorem 2** *Given an initial set-based abstraction \mathcal{N} and assuming $\mathcal{N}(g_m)$ satisfies the Conditions (1), (2),
 303 and (3) of Definition 3, there exists a finite number of refinements under which a reachability-aware abstraction
 304 can be computed.*

305 Theorem 2 establishes the effectiveness of the representation refinement process that gradually builds a
 306 reachability-aware abstraction. The theorem follows from Lemma 1. In practice, the assumption that $\mathcal{N}(g_m)$
 307 verifies criteria 1,2,4 of Def.3 is reasonable since the goal g^* is known in M . That is to say, $\mathcal{N}(g_m)$ could
 308 correspond to a region whose center is g^* and radius is R_{max} .