
Direct then Diffuse: Incremental Unsupervised Skill Discovery for State Covering and Goal Reaching

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Learning meaningful behaviors in the absence of a task-specific reward function
2 is a challenging problem in reinforcement learning. A desirable unsupervised
3 objective is to learn a set of diverse skills that provide a thorough coverage of
4 the state space while being directed, i.e., reliably reaching distinct regions of the
5 environment. At test time, an agent could then leverage these skills to solve sparse
6 reward problems by performing efficient exploration and finding an effective goal-
7 directed policy with little-to-no additional learning. Unfortunately, it is challenging
8 to learn skills with such properties, as diffusing (e.g., stochastic policies performing
9 good coverage) skills are not reliable in targeting specific states, whereas directed
10 (e.g., goal-based policies) skills provide limited coverage. In this paper, inspired
11 by the mutual information framework, we propose a novel algorithm designed
12 to maximize coverage while ensuring a constraint on the directedness of each
13 skill. In particular, we design skills with a decoupled policy structure, with a first
14 part trained to be directed and a second diffusing part that ensures local coverage.
15 Furthermore, we leverage the directedness constraint to adaptively add or remove
16 skills as well as incrementally compose them along a tree that is grown to achieve a
17 thorough coverage of the environment. We illustrate how our learned skills enables
18 to efficiently solve sparse-reward downstream tasks in navigation and continuous
19 control environments, where it compares favorably with existing baselines.

20 1 Introduction

21 Deep reinforcement learning (RL) algorithms have been shown to effectively solve a wide variety of
22 complex problems [e.g., 30, 6, 39, 16, 2, 36]. However, they are often designed to solve one single
23 task at a time and they need to restart the learning process from scratch for any new problem, even
24 when it is defined on the very same environment (e.g., navigating to different locations in the same
25 apartment). Recently, unsupervised RL (URL) has been proposed as an approach to address this
26 limitation. In URL, the agent first interacts with the environment without any extrinsic reward signal.
27 Afterward, the agent leverages the experience accumulated during the unsupervised learning phase to
28 efficiently solve a variety of downstream tasks defined on the same environment.

29 In this paper, we consider the URL setting where the agent starts from an initial state s_0 and it resets
30 to it every time the policy terminates. We focus on sparse-reward downstream tasks, which require
31 effective exploration (i.e., via a thorough coverage of the state space) to find the goal as well as
32 learning a policy reliably reaching the goal (i.e., a directed policy).

33 We build on the insight that *mutual information* (MI) effectively formalizes the dual objective of
34 learning skills that both cover and navigate the environment efficiently [e.g., 15]. Specifically, given
35 the state variable S and some variables Z on which the skill policies are conditioned, MI is defined as

$$\mathcal{I}(S; Z) = \underbrace{\mathcal{H}(S)}_{\text{coverage}} - \underbrace{\mathcal{H}(S|Z)}_{\text{directedness}} = \mathcal{H}(Z) - \mathcal{H}(Z|S), \quad (1)$$

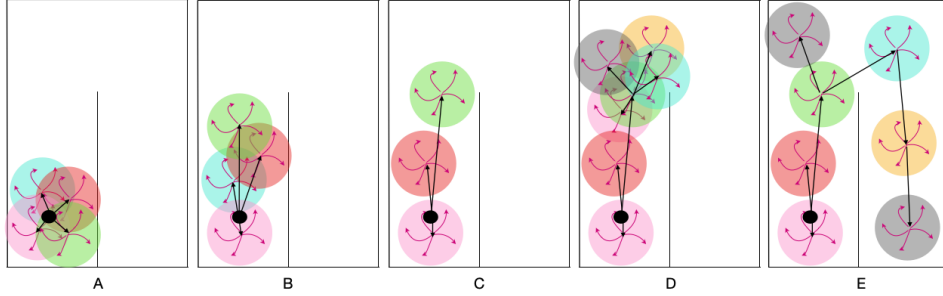


Figure 1: Overview of UPSIDE. The black dot corresponds to the initial state s_0 . (A) A set of random skills is initialized, each skill being composed of a *directed* part (illustrated as a black arrow) and a *diffusing* part (red arrows), which induces a local coverage (colored circles). (B) The policies associated to the directed part of each skill are then updated to maximize the discriminability of the states reached by their diffusing part (Sect. 3.1). (C) The least discriminable skills are iteratively removed while the policies of the remaining skills are re-optimized. This is executed until the discriminability of each skill satisfies a given constraint (see Sect. 3.2). In this example three skills are kept. (D) One of these learned skill is then used as basis to add new skills, which are then optimized following the same procedure. For the “red” and “purple” skills, UPSIDE is not able to find sub-skills of sufficient quality and thus they are not expanded any further. (E) At the end of the process, UPSIDE has created a tree of directed skills covering the state space (Sect. 3.3). These covering skills can then be used to solve downstream tasks. Moreover, the discriminator learned together with the skills can be used to select the skill to reach any specific goal region, where the directed parts get close to the goal, while the diffusing part provides the local coverage to attain the goal. The complete algorithm is detailed in Sect. 3.4 and Appendix.

36 where \mathcal{I} denotes the MI and \mathcal{H} is the entropy function. The first expression, known as the forward
 37 form of MI, explicitly balances the two sought-after properties of *coverage* — captured by the entropy
 38 over the state space $\mathcal{H}(S)$ — and *directedness*, i.e., the ability to reach specific states S depending
 39 on Z — captured by the negative conditional entropy $-\mathcal{H}(S|Z)$. The second expression of (1), often
 40 easier to optimize and referred to as the reverse form, stipulates that the skills should be sampled as
 41 diversely as possible while being discriminable.

42 Maximizing (1) has been shown to be a powerful approach for encouraging exploration in RL
 43 [20, 32] and for unsupervised skill discovery [e.g., 15, 12, 1, 38, 10]. Nonetheless, learning skills
 44 that maximize the MI is a challenging optimization problem. Several approximations have been
 45 proposed to simplify the problem at the cost of possibly deviating from the original objective of
 46 coverage and directedness (see Sect. 4 for a review of related work). In this paper, we propose
 47 UPSIDE (*UnsuPervised Skills that direct then Diffuse*) to learn skills that can be effectively used to
 48 solve goal-based downstream tasks. Our solution builds on the following components (see Fig. 1 for
 49 an illustration of UPSIDE):

- 50 • *Skill structure.* In order to balance coverage and directedness, we design skills composed of two
 51 parts: **1**) a *directed* part that is trained to reach a distinct region of the environment, and **2**) a
 52 *diffusing* part that covers the states around the region attained by the first part.
- 53 • *Optimization.* We further strengthen the coverage and directedness properties of the skills by
 54 turning the MI objective into a constrained optimization problem designed to maximize coverage
 55 under the constraint that *each* skill achieves a minimum level of discriminability. This in turn
 56 enables UPSIDE to adaptively add skills to improve coverage, when all the initial skills meet the
 57 constraint, or remove those that violate the constraint to guarantee that each skill is directed and
 58 reaches a distinct region of the environment.
- 59 • *Tree structure.* When the agent starts from a fixed initial state, the skills’ length is a crucial
 60 parameter, where short skills do not allow for proper coverage, and long skills are difficult to train.
 61 In UPSIDE we consider short skills to make the optimization easier, while composing them along a
 62 tree structure that ensures an adaptive and deep coverage of the environment.

63 We study how our learned skill structure enables to both perform efficient exploration and learn
 64 effective goal-reaching policies in a variety of navigation and continuous control environments
 65 (including MuJoCo’s reacher) and we compare its performance to relevant baselines.

66 2 Setting

67 We consider the URL setting where the agent interacts with a Markov decision process (MDP) M
 68 with state space \mathcal{S} , action space \mathcal{A} , dynamics $p(s'|s, a)$, and **no reward**. The agent starts each
 69 episode from a designated initial state $s_0 \in \mathcal{S}$. Upon termination of the chosen policy, the agent is
 70 then reset to s_0 . This setting is particularly challenging from an exploration point of view since the
 71 agent cannot rely on the initial distribution to cover the state space.

72 We recall the MI-based unsupervised skill discovery approach [see e.g., 15]. Denote by Z some
 73 (latent) variables on which the skills of length T are conditioned. There are three optimization
 74 variables: (i) the support of the skills denoted by $|Z|$ (we consider it to be discrete so $|Z|$ is the
 75 number of skills), (ii) the policy $\pi(z)$ associated to skill z , and (iii) the sampling rule ρ (i.e., $\rho(z)$
 76 is the probability of sampling skill z at the beginning of the episode). Let the variable S_T be the
 77 random (final) state induced by sampling a skill z from ρ and executing the associated policy $\pi(z)$
 78 from s_0 for an episode. We denote by $p_{\pi(z)}(s_T)$ the distribution over (final) states induced by
 79 executing the policy of skill z , by $p(z|s_T)$ the probability of z being the skill to induce state s_T , and
 80 let $\bar{p}(s_T) = \sum_{z \in Z} \rho(z) p_{\pi(z)}(s_T)$. Then maximizing the MI between Z and S_T can be written as

$$\begin{aligned} \max_{|Z|, \rho, \pi} \mathcal{I}(S_T; Z) &= \mathcal{H}(S_T) - \mathcal{H}(S_T|Z) = - \sum_{s_T} \bar{p}(s_T) \log \bar{p}(s_T) + \sum_{z \in Z} \rho(z) \mathbb{E}_{s_T} [\log p_{\pi(z)}(s_T)] \\ &= \mathcal{H}(Z) - \mathcal{H}(Z|S_T) = - \sum_{z \in Z} \rho(z) \log \rho(z) + \sum_{z \in Z} \rho(z) \mathbb{E}_{s_T} [\log p(z|s_T)], \quad (1) \end{aligned}$$

81 where in the expectations $s_T \sim p_{\pi(z)}(s_T)$. As discussed in Sect. 1, learning the optimal $|Z|$, ρ , and π
 82 is a challenging problem [see e.g., 15, 12, 10].

83 3 Algorithm Structure

84 UPSIDE is based on three main components: **a**) the skill learning corresponding to stage A and B of
 85 Fig. 1 and described in Sect. 3.1, **b**) a constrained optimization problem used to optimize the number
 86 of skills (stage C and Sect. 3.2) and **c**) a tree-building procedure (stage D and Sect. 3.3). Together,
 87 these components allow UPSIDE to discover skills that combine coverage and directedness.

88 3.1 Skill Structure and Optimization

89 As shown in e.g., [12, 38, 46], the level of stochasticity of each skill (e.g., induced via a regularization
 90 on the entropy over the actions) plays a key role in trading off coverage and directedness. In fact,
 91 while randomness promotes broader coverage, it may compromise the directedness of the skills.
 92 In fact, a highly stochastic skill tends to induce a distribution $p_{\pi(z)}(s_T)$ over final states with high
 93 entropy (thus decreasing $-\mathcal{H}(S_T|Z)$), which prevents the skill to be reusable in solving sparse-reward
 94 downstream tasks where the objective is to reliably reach a specific goal state of the environment.
 95 Determining *how much* stochasticity to inject to adequately balance both objectives and optimize (1)
 96 is a difficult problem.¹

97 We propose to design skills with a *decoupled policy structure*:

- A *directed* part (of length T) with low stochasticity and trained to reach a specific region of the environment. It is responsible for increasing the $-\mathcal{H}(S|Z)$ term in (1).
- A *diffusing* part (of length H) with high stochasticity to promote local coverage of the states around the region reached by the directed part. It is responsible for increasing the $\mathcal{H}(S)$ term in (1).

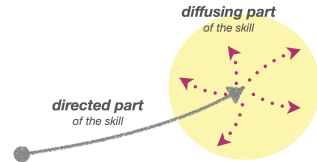


Figure 2: Directed and diffusing parts of the skill.

99 Similar to prior work [e.g., 15, 12], the policy associated to the directed part of skill z is trained to max-
 100 imize an intrinsic reward $r_z(s) \approx p(z|s)$,² where $p(z|s)$ measures the “discriminability” of the skill z
 101 given the state s . More formally, $\pi(z)$ maximizes the cumulative reward $\mathbb{E}_{\pi(z)} [\sum_{t=T+1}^{T+H} r_z(s_t)]$
 102 over the states traversed by the policy during the diffusing part. In practice, we also add a small
 103 entropy regularization $\mathcal{H}(\pi(\cdot|z, s_t))$ to the directed policy in order to ensure a minimum level of

¹In RL, stochasticity is injected at “train time” to boost *exploration* or improve *robustness*, while the policy executed at “test time” is deterministic. Here we refer to stochasticity introduced to better optimize (1).

²Although [15, 12] employ rewards in the log domain, we find that using a reward that is a non-linear transformation into $[0, 1]$ works better in practice, as also observed in [42, 5]. Furthermore, in practice we replace $p(z|s)$ by the predictions of a learned discriminator $q_\phi(z|s)$ as explained in Sect. 3.4.

104 exploration and make the learning more robust. For the diffusing part, we rely on a simple random
 105 walk policy (i.e., a stochastic policy with uniform distribution over actions).

106 Intuitively, the diffusing part defines a cluster of states that is used as a goal for the directed part.
 107 This allows us to “ground” the latent variable representations of the skills Z to specific regions of
 108 the environment (i.e., the clusters). As a result, maximizing the MI over such skills can be seen as
 109 learning a set of “cluster-conditioned”, and thus directed, policies.

110 3.2 Skill Support and Sampling Rule

111 The MI objective (1) crucially depends on the number of skills ($|Z|$) and the distribution $\rho(z)$.
 112 Unfortunately, it is been shown [e.g., 10] that solving (1) is particularly challenging. In order to
 113 simplify the optimization and the associated learning problem, we modify (1) in two ways.

114 First, coherently with the skill optimization detailed in Sect. 3.1, the random variable S in the
 115 conditional entropy is any state reached during the diffusing part of the skill and not just the terminal
 116 state. More formally, we denote by S_{diff} the random variable and its distribution for a specific skill z
 117 is $p_{\pi(z)}(s_{\text{diff}}) = 1/H \sum_{t=T+1}^{T+H} p_{\pi(z)}(s_t)$, i.e., the distribution over states obtained by averaging the
 118 distributions at any of the steps in the diffusing part. Similarly, $p(z|s_{\text{diff}})$ now denotes the probability
 119 of z being the skill to traverse s_{diff} during its diffusing part. As a result, training the skills to maximize
 120 MI naturally leads the diffusing parts to “push” the directed parts away so as to reach diverse regions
 121 of the environment. The combination of “global” coverage of the directed parts and “local” coverage
 122 of the diffusing part ensures that the whole environment is properly visited with $|Z| \ll S$ skills.³

123 Second, we introduce an alternative problem that simplifies the optimization while preserving the
 124 coverage and directedness properties of MI. This is achieved by introducing a stronger requirement
 125 on the discriminability. While the conditional entropy term $-\mathcal{H}(Z|S)$ in (1) promotes the discrim-
 126 inability of skills *on average*, we argue that a more suitable objective is to *constrain* each skill to
 127 achieve a *minimum* level of discriminability. First, we move from the average to the minimum over
 128 skills by lower bounding the conditional entropy as

$$-\mathcal{H}(Z|S_{\text{diff}}) = \sum_{z \in Z} \rho(z) \mathbb{E}_{s_{\text{diff}}} [\log p(z|s_{\text{diff}})] \geq \min_{z \in Z} \mathbb{E}_{s_{\text{diff}}} [\log p(z|s_{\text{diff}})], \quad (2)$$

129 which leads to the following optimization (assuming π is fixed for convenience)

$$\max_{|Z|=N, \rho} \left\{ \mathcal{H}(Z) + \min_{z \in [N]} \mathbb{E}_{s_{\text{diff}}} [\log p(z|s_{\text{diff}})] \right\}, \quad (3)$$

130 where with an abuse of notation we use $z \in [N]$ to denote all skills in a set Z with cardinality N .
 131 Since (3) is a lower bound to MI, it tends to promote the same type of covering and directed skills.
 132 Furthermore, (2) no longer depends on the distribution over skills and the entropy term $\mathcal{H}(Z)$ is
 133 maximized by setting ρ to the uniform distribution over N skills (i.e., $\max_{\rho} \mathcal{H}(Z) = \log(N)$), thus
 134 simplifying the optimization, which now only depends on N .

135 While optimizing (3) promotes a cardinality N such that all skills have good discriminability, a more
 136 convenient formulation is to explicitly set a minimum level of discriminability for all skills through
 137 the following constrained optimization problem:

$$\max_{N \geq 1} \log(N) \quad \text{s.t.} \quad \min_{z \in [N]} \mathbb{E}_{s_{\text{diff}}} [\log p(z|s_{\text{diff}})] \geq \log \eta. \quad (4)$$

138 where η is a parameter that defines the discriminability threshold. A skill z is said to be η -*consolidated*
 139 if it satisfies the constraint. Crucially, let $P_N := \min_{z \in [N]} \mathbb{E}_{s_{\text{diff}}} [\log p(z|s_{\text{diff}})]$, then the sequence
 140 $(P_N)_{N \geq 1}$ is non-increasing with $P_1 = 0$ (i.e., the more skills the harder it is to meet the constraint).
 141 As a result, (4) can be optimized following a simple greedy strategy incrementally adding skills until
 142 the constraint is violated. The optimal N thus defines the *effective number* of η -consolidated skills and
 143 it corresponds to the largest number of skills that is guaranteed to display sufficient discriminability.
 144 Alternatively, we can interpret (4) as finding the largest number of clusters (i.e., the region reached
 145 by the directed part of a skill and covered by its associated diffusing part) with a minimum level of
 146 inter-cluster distance. This effect is qualitatively illustrated in Fig. 1, where the states attained by the
 147 directed part of the skills attain different regions that are locally covered by their diffusing parts.

³Notice that (1) is maximized by setting $|Z| = |S|$ (since $\max_Y \mathcal{I}(X, Y) = \mathcal{I}(X, X) = \mathcal{H}(X)$), i.e., where each skill is a goal-conditioned policy reaching a different state. This implies having as many policies as states, which makes the learning particularly challenging as the complexity of the environment increases.

Algorithm 1: UPSIDE

Initialize: Discriminability threshold $\eta \in (0, 1)$, branching factor $N_0 \geq 1$, patience K

Initialize: Tree \mathcal{T} initialized as a root node indexed by 0, queue of parent nodes $\mathcal{W} = \{0\}$.

```
while  $\mathcal{W} \neq \emptyset$  do // tree expansion
1  Dequeue a node/skill  $w \in \mathcal{W}$  and expand  $\mathcal{T}$  at  $w$  by adding a set  $\mathcal{C}(w)$  of  $N_0$  nodes/skills
2  Create random policies  $\pi_z, \forall z \in \mathcal{C}(w)$ 
3  Initialize discriminator  $q_\phi$  with  $|\mathcal{T}|$  classes
4  Continue = true; Saturated = false
5  while Continue do
6    for  $K$  iterations do
7      Sample a skill  $z$  from  $\mathcal{T}$  at random
8      Extract the sequence of nodes  $z_{(1)}, \dots, z$  in  $\mathcal{T}$  leading to  $z$ 
9      Execute the composed (directed part) policy  $(\pi_{z_{(1)}}, \dots, \pi_z)$  followed by the diffusing part
10     Add states observed during the diffusion part to state buffer  $\mathcal{B}_z$ 
11     Update discriminator  $q_\phi$  with SGD on  $\mathcal{B}_z$  to predict label  $z$ 
12     if  $z \in \mathcal{C}(w)$  then // Update only new policies, other policies kept fixed
13       | Update policy  $\pi_z$  using SAC to optimize the discriminator reward as in Sect. 3.1.
14     Compute the skill-discriminability  $d(z) = \hat{q}_\phi^{(B)}(z) = \frac{1}{|\mathcal{B}_z|} \sum_{s \in \mathcal{B}_z} q_\phi(z|s)$  for all  $z \in \mathcal{C}(w)$ 
15     if  $\min_{z \in \mathcal{C}(w)} d(z) < \eta$  then // Node removal
16       | Remove the node/skill  $z = \arg \min_{z \in \mathcal{C}(w)} d(z)$  from  $\mathcal{C}(w)$  and  $\mathcal{T}$ 
17       | Set Saturated = true
18     else if not Saturated then
19       | Add one new node/skill to  $\mathcal{C}(w)$  and  $\mathcal{T}$ 
20     else
21       | Set Continue = false
22  Enqueue in  $\mathcal{W}$  the consolidated nodes  $\mathcal{C}(w)$ 
```

148 3.3 Composing Skills in a Tree Structure

149 The MI optimization problem as well as our constrained variant (4) depend on the initial state s_0
150 and on the length of each skill. Although these quantities are usually predefined and only appear
151 implicitly in the equations, they have a crucial impact on the obtained behavior. In fact, resetting after
152 each skill execution unavoidably restricts the coverage to a radius of at most $T + H$ steps around s_0 .
153 This may suggest to set T and H to a large value. However, increasing the horizon makes the training
154 of the skills more challenging, as learning π would require solving a difficult RL problem itself.

155 Instead, we propose to “extend” the length of the skills through composition. Indeed, the decoupled
156 skill structure and the constraint in (4) entail that the directed part of each of the η -consolidated skills
157 reliably reach a specific (and distinct) region of the environment and it is thus re-usable and amenable
158 to composition. We propose to chain the directed part of the skills in order to reach further and further
159 parts of the state space. Specifically, we build a growing tree, where the root is the initial state s_0 , the
160 edges represent the directed part of the skills, and the nodes represent the diffusing part of skills. As
161 such, whenever a skill z is selected, the directed part of all the policies associated to its predecessor
162 skills in the tree are executed first (see Fig. 1 for an illustration of the tree structure).

163 As a result, the agent naturally builds a curriculum on the episode lengths, which grow as the sequence
164 $(iT + H)_{i \geq 1}$. As such, it does not require prior knowledge on an adequate horizon of the downstream
165 goal-based task.⁴ Here this knowledge is replaced by T and H which are more environment-agnostic
166 and task-agnostic quantities, as their choice rather has an impact on the size and shape of the learned
167 tree (e.g., the smaller T and H the bigger the tree).

168 3.4 The UPSIDE Algorithm

169 We are now ready to introduce UPSIDE, which provides a specific implementation of the components
170 described before (see Fig. 1 for a qualitative illustration and Algorithm 2 for the detailed pseudo-code).

171 We perform standard approximations to make the constraint in (4) easier to estimate. We approximate
172 the unknown posterior $p(z|s)$ with a learned discriminator $q_\phi(z|s)$ with parameters ϕ . We also

⁴See e.g., the discussion in [33] on the “importance of properly choosing the training horizon in accordance with the downstream-task horizon the policy will eventually face.”

173 remove the logarithm from the constraint to have an estimation range of $[0, 1]$ and thus lower
 174 variance². Finally, we replace the expectation over s with an empirical estimate $\widehat{q}_\phi^{(B)}(z)$ averaging the
 175 value of the discriminator evaluated on the last B states observed while executing the diffusing part
 176 of z . Integrating these approximations in (4) leads to

$$\max_{N \geq 1, \pi} N \quad \text{s.t.} \quad \min_{z \in [N]} \widehat{q}_\phi^{(B)}(z) \geq \eta. \quad (5)$$

177 As discussed in Sect. 3.2, this problem can be conveniently optimized using a greedy strategy. We
 178 then integrate the optimization of (5) into an adaptive tree expansion strategy: **(Generating new**
 179 **skills)** Given a tree structure as described in Sect. 3.3, we expand the tree at a leaf w by adding N_0
 180 new nodes/skills following a breadth-first-search approach (lines 1, 2). Then **(Skill Learning)** the
 181 new skills are optimized by: **i)** sampling random skills in the tree to update the discriminator (lines
 182 7-11), and **ii)** by updating the policies to optimize the discriminability reward (Sect. 3.1) computed
 183 using the discriminator (lines 13). To speed-up convergence, we only update the policies that have be
 184 added to the tree structure, keeping all the previous policies fixed (line 12). Note that in the update of
 185 the discriminator we leverage the states observed in previous phases of the algorithm by maintaining
 186 a (small) replay buffer of states for each skill. **(Node Consolidation)** After a *patience* period (line 6),
 187 if all skills are η -consolidated, we tentatively add more skills to the leaf w (line 18). On the other
 188 hand, if any skill does not meet the discriminability threshold, we remove it and consolidate the
 189 remaining skills into the tree (lines 16, 17) and we repeat the process.

190 **Model selection.** A core aspect of any RL algorithm is *model selection*, i.e., finding the best
 191 configuration of hyperparameters. In URL with no prior knowledge of the downstream task(s), it
 192 is non-trivial to devise an adequate criterion for model selection and this aspect is rarely addressed,
 193 despite being crucial in practice. For instance, while the coverage of the state space may be a good
 194 proxy for the performance of a URL algorithm [see e.g., 10], it may be difficult to measure in
 195 continuous problems. Interestingly, our optimization problem directly provides a single, task-agnostic
 196 and environment-agnostic criterion for model selection, which is the number N of η -consolidated
 197 skills discovered by the agent. Indeed in all of our experiments we simply select the model (i.e., set
 198 of hyperparameters) that maximizes N . This is a significant advantage w.r.t. existing methods, such
 199 as VIC and DIAYN, for which no principled approach to model selection is provided.

200 4 Related work

201 Unsupervised Reinforcement Learning methods can be broadly decomposed according to the way
 202 they summarize the experience accumulated during the unsupervised phase into reusable knowledge
 203 to solve downstream tasks. This includes both off-policy model-free [e.g., 34] and model-based
 204 [e.g., 37] methods that seek to populate a representative replay buffer and build accurate value or
 205 model estimates, that are used to solve a given downstream task in a zero- or few-shot manner.
 206 The accumulated experience during train time can also be compressed into a low-dimensional
 207 representation for value functions as well as policies and to improve exploration [e.g., 45]. An
 208 alternative line of work focuses on the discovery of a set of skills in an unsupervised manner. Our
 209 approach falls in this category, on which we now focus our related work review.

210 Skill discovery based on MI maximization was first proposed in VIC [15], where only the final states
 211 of each trajectory are considered in the reverse form of (1) and where both the skills and their sampling
 212 rules are simultaneously learned (with a fixed support $|Z|$, i.e., a fixed number of skills). DIAYN [12]
 213 fixes the sampling rule to be uniform, and weighs the skills with an action-entropy coefficient (i.e., it
 214 additionally minimizes the MI between actions and skills given the state), so as to push the skills
 215 away from each other and enhance coverage. DADS [38] learns skills that are not only diverse but
 216 also predictable by learned dynamics models, by using a generative model over observations (rather
 217 than over skills) and optimizing a forward form of MI, namely $\mathcal{I}(s'; z|s)$ between the next state s'
 218 and current skill z (with continuous latent) conditioned on the current state s . EDL [10] shows that
 219 existing skill discovery approaches can provide insufficient coverage, and instead proposes to rely on
 220 a fixed distribution over states $p(s)$ which is either provided by an oracle or learned. In SMM [24], the
 221 MI formalism is used to learn a policy for which the state marginal distribution matches a given target
 222 state distribution (e.g., uniform), which can be seen as a more scalable way of tackling the problem of
 223 maximum entropy over the state space [19], and as a way to encourage skills to go through unknown
 224 state regions. Other MI-based skill discovery methods include [13, 18, 31, 5, 43], as well as [44, 27]
 225 which investigate skill discovery in non-episodic settings.

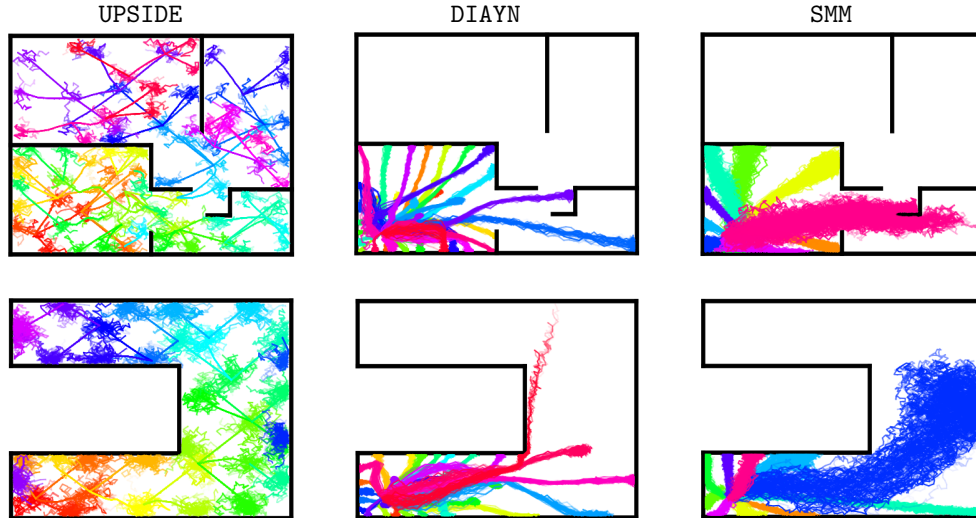


Figure 3: UPSIDE, DIAYN-curriculum and SMM-10 skills learned in a bottleneck maze (*Top*) and a U-maze (*Bottom*). For both DIAYN and SMM we report the stochastic execution of the learned skills and for UPSIDE we report the deterministic directed parts (that are composed) followed by the (stochastic) diffusing part, which is the same protocol used to evaluate coverage.

226 Our approach shares a similar motivation to prior MI-based works of targeting skills that are both
 227 directed and state-covering. In particular, the decoupled structure introduced in Sect. 3.1 can be seen
 228 as a more suitable way to achieve the objective of improving the coverage of VIC as done in DIAYN
 229 and SMM, without compromising the directedness of the skills.

230 While most skill discovery approaches consider a fixed number of skills, a curriculum with increasing
 231 number of skills is studied in [1, 3]. Our discriminability constraint is what enables skills to be
 232 composed along a tree structure, which allows increases or decreases the support of available skills
 233 depending on the region of the state space.

234 Recently, [46] proposed a hierarchical RL method that discovers abstract and task-agnostic skills
 235 while jointly learning a higher-level policy which is trained to maximize environment reward. Our
 236 approach builds on a similar promise of composing skills instead of resetting to s_0 after each execution,
 237 yet we articulate the composition differently, by exploiting the direct-then-diffuse structure to ground
 238 learned skills to the state space instead of being abstract.

239 In addition, approaches such as DISCERN [42] and Skew-Fit [34] learn a goal-conditioned policy in
 240 an unsupervised way with an MI objective. As explained in [10, Sect. 5], this can be interpreted as a
 241 skill discovery approach with latent $Z = S$, i.e., where each goal state can define a different skill.
 242 Conditioning on either goal states or abstract latent skills forms two extremes of the spectrum of
 243 unsupervised RL. We target an intermediate approach, seeking to benefit from the groundedness of
 244 the latent skill Z and the states S (and thus amenability to composition) of goal-conditioned RL, and
 245 from the reduced search space and sampling ease of skill-based RL.

246 An alternative approach to skill discovery builds on “spectral” properties of the dynamics of the
 247 environment. This includes eigenoptions [28, 29] and covering options [22, 23], as well as the
 248 algorithm of [4] that builds a discrete graph representation which learns and composes spectral skills.

249 5 Experiments

250 In this section, we investigate the following questions: **i)** Can the adaptive tree structure of UPSIDE in-
 251 crementally cover an unknown environment while preserving directedness of the skills? **ii)** Following
 252 the unsupervised phase, how can UPSIDE be leveraged to solve goal-based downstream tasks?

253 We report results on: **a)** Navigation problems in continuous mazes, where actions represent the desired
 254 shift in x and y coordinates; **b)** A difficult instance of CartPole, where the cart starts with zero speed
 255 and the pole is oriented downside; **c)** The Reacher [41] problem using the MuJoCo implementation
 256 in Gym [8]. In all environments, the per-dimension action space is in $[-1; +1]$.

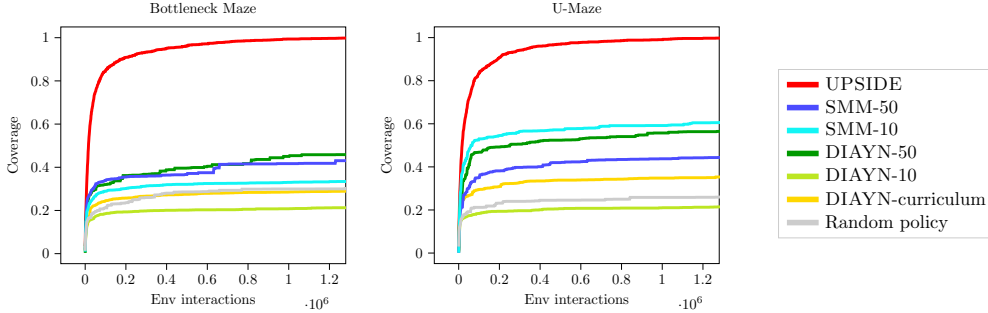


Figure 4: Normalized coverage in U-maze and bottleneck.

257 We compare to different baselines. DIAYN- K , where K is a fixed number of skills, is the original
 258 algorithm proposed in [12]. DIAYN-Curriculum is a variant where the number of skills is automatically
 259 tuned following the same procedure as in UPSIDE ensuring a good discriminability. We also compare
 260 to SMM [24], which is similar to DIAYN, but it includes an exploration bonus encouraging the policies
 261 to visit rarely encountered states. In our implementation, the exploration bonus is obtained by
 262 maintaining a multinomial distribution over “buckets of states” obtained by discretization, resulting
 263 in an computation-efficient and stable implementation that is more stable than the original VAE-based
 264 method. UPSIDE and all baselines are implemented with Soft-Actor Critic (SAC) [17].

265 **Unsupervised Phase.** We run all methods until convergence. We then do model selection according
 266 to the criterion of either the final number of skills for UPSIDE and DIAYN-curriculum and the final
 267 average discriminability for DIAYN- K and SMM. To compute the coverage, we perform rollouts by
 268 first sampling a skill uniformly at random and executing its associated policy until termination. We
 269 discretize states into buckets (50 interval per dimension for mazes and 10 for control environments)
 270 and report the proportion of buckets reached by each method as a function of the total number of
 271 steps executed in the environment over multiple rollouts. Since only a small portion of the discretized
 272 states can be reached, we normalize the coverage such that the best method obtains 1.

273 We consider two topologies of mazes with size (height and width) 50 such that exploration is non-
 274 trivial (i.e., a random policy is only able to cover a small part of the state space): a U-shaped maze
 275 and a Bottleneck maze (which is a harder version of the one in [10, Fig. 1] which is only of size 10
 276 for the same action space). In Fig. 3 we show that UPSIDE succeeds in covering the near-entirety
 277 of the state space by creating a tree of directed skills. Moreover, UPSIDE created directed skills
 278 with a low entropy, while the two baselines tend to create skills that are more stochastic. This is
 279 particularly evident for SMM, due to the state-entropy exploration bonus, that while it encourages
 280 broader coverage makes skills less directed.

281 In Fig. 4 we report the coverage on the Bottleneck maze and U-Maze. For UPSIDE, executing a
 282 skill corresponds to executing the directed part of all the “parent” skills in the tree and concluding
 283 with the diffusion part of the skill. SMM achieves better coverage than DIAYN thanks to the increased
 284 level of stochasticity (diffusion) of its skills. UPSIDE outperforms both by reaching regions of the
 285 environment that are not be achieved by other methods. Here, we plot UPSIDE with $T = 10$ and
 286 $H = 10$, but we found UPSIDE to be robust to these parameters as shown in the supplementary.

287 Results are similar in the CartPole problem (see Fig. 5) where UPSIDE (with $T = 20$ and $H = 20$)
 288 obtains better coverage than baselines. On the other hand, in Reacher (see Fig. 5), DIAYN-50
 289 outperforms UPSIDE in terms of coverage. This can be explained by the fact that, in this environment,
 290 highly stochastic skills provide a good coverage. Nonetheless, this comes at the cost of very low
 291 discriminability (rightmost plot), which suggests DIAYN-50 skills have poor directedness. On the
 292 other hand, UPSIDE (and DIAYN-curriculum) achieves much larger discriminability by removing
 293 redundant skills and favoring more directed policies.

294 **Downstream Tasks.** Following the unsupervised phase, UPSIDE has learned a tree of skills. We
 295 now investigate how these skills are used to tackle a downstream task. In that setting, we propose to
 296 use skill-based approaches (i.e UPSIDE, DIAYN and SMM) in the following way: a) (exploration) first
 297 we sample rollouts over the different skills. b) We then select the best skill based on the maximum
 298 cumulative reward collected and c) we fine-tune this skill to maximize the reward. We report results
 299 on mazes (additional results are provided in the supplementary). We consider a sparse positive reward

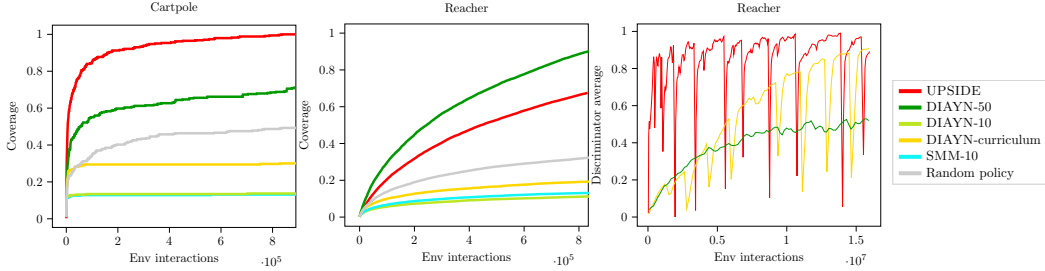
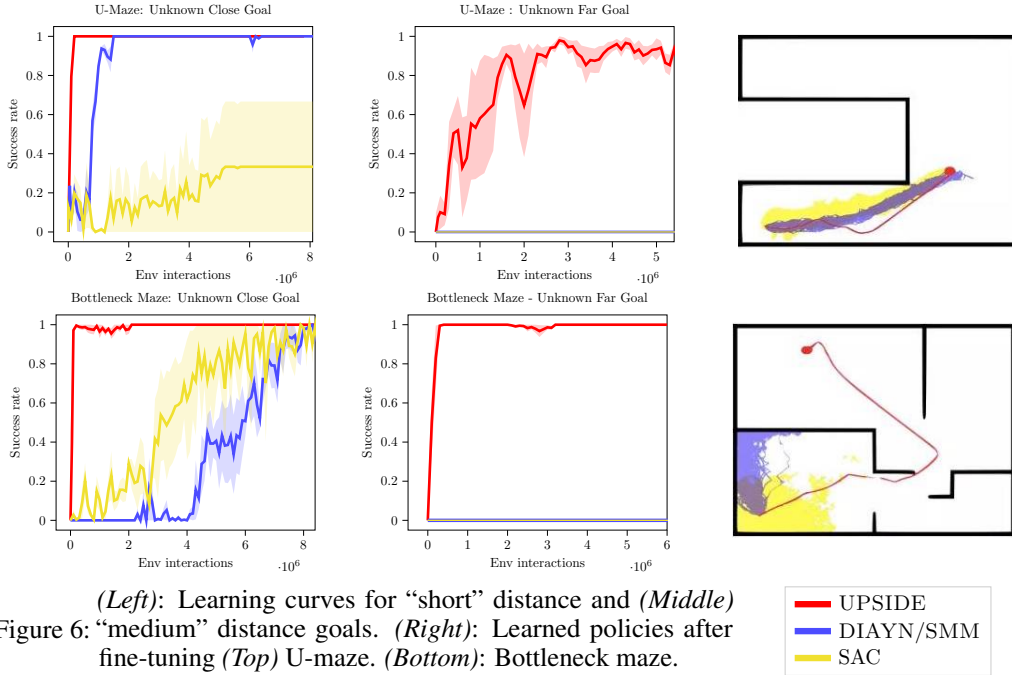


Figure 5: Normalized coverage in Cartpole (*Left*) and Reacher (*Middle*). (*Right*) Average discriminability of the skills during training in Reacher.

300 when reaching a particular defined goal.⁵ We consider goals at different distances from the initial
 301 state s_0 , the further, the harder. Fig. 6 shows the learning curves obtained when fine-tuning the best
 302 skill for the different models and compare to a classical SAC algorithm where a single policy is
 303 learned from scratch. DIAYN/SMM means we use the best state-covering policies between DIAYN and
 304 SMM. For the “close” goal setting, both UPSIDE and DIAYN/SMM are able to learn to reach this goal
 305 efficiently while SAC solves the task only for some of the training runs. Note that we do not show
 306 DIAYN performance since it is lower than the SMM one. For the “far” goal setting, only UPSIDE learns
 307 to reach this goal. Obtained trajectories are illustrated in Fig. 6.



(*Left*): Learning curves for “short” distance and (*Middle*) “medium” distance goals. (*Right*): Learned policies after fine-tuning (*Top*) U-maze. (*Bottom*): Bottleneck maze.

308 6 Conclusion

309 We introduced UPSIDE, a novel algorithm for unsupervised skill discovery designed to trade off
 310 between coverage and directedness and develop a tree of skills that can be used to both perform
 311 efficient exploration of the environment and learn effective goal-directed policies. Natural venues for
 312 future investigation are: **1)** The diffusing part of each skill could be explicitly trained to maximize
 313 local coverage; **2)** UPSIDE assumes a good representation of the state is provided as input, it would
 314 be interesting to pair UPSIDE with effective representation learning techniques to tackle problems
 315 with high-dimensional input (e.g., image-based RL); **3)** While UPSIDE is grounded on the solid
 316 principle of MI maximization, a more thorough theoretical investigation is needed to explicitly link
 317 the optimization problem and its approximations to the downstream performance.

⁵Notice that if the goal was known, the learned discriminator could be directly used to identify the most promising skill to fine-tune.

318 **References**

- 319 [1] J. Achiam, H. Edwards, D. Amodei, and P. Abbeel. Variational option discovery algorithms.
320 *arXiv preprint arXiv:1807.10299*, 2018.
- 321 [2] M. Andrychowicz, D. Crow, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin,
322 P. Abbeel, and W. Zaremba. Hindsight experience replay. In *NIPS*, 2017.
- 323 [3] A. Aubret, L. Matignon, and S. Hassas. Elsim: End-to-end learning of reusable skills through
324 intrinsic motivation. *arXiv preprint arXiv:2006.12903*, 2020.
- 325 [4] A. Bagaria, J. Crowley, J. W. N. Lim, and G. Konidaris. Skill discovery for exploration and
326 planning using deep skill graphs. 2021.
- 327 [5] K. Baumli, D. Warde-Farley, S. Hansen, and V. Mnih. Relative variational intrinsic control.
328 *arXiv preprint arXiv:2012.07827*, 2020.
- 329 [6] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An
330 evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279,
331 2013.
- 332 [7] D. Bertsekas. *Dynamic programming and optimal control*, volume 2. 1995.
- 333 [8] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba.
334 Openai gym, 2016.
- 335 [9] A. Z. Broder and A. R. Karlin. Bounds on the cover time. *Journal of Theoretical Probability*,
336 2(1):101–120, 1989.
- 337 [10] V. Campos, A. Trott, C. Xiong, R. Socher, X. Giro-i Nieto, and J. Torres. Explore, discover and
338 learn: Unsupervised discovery of state-covering skills. In *International Conference on Machine*
339 *Learning*, 2020.
- 340 [11] M. Diaz, L. Paull, and P. S. Castro. Loco: Adaptive exploration in reinforcement learning via
341 local estimation of contraction coefficients. *ICLR Workshop SSL-RL*, 2021.
- 342 [12] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine. Diversity is all you need: Learning skills
343 without a reward function. In *International Conference on Learning Representations*, 2019.
- 344 [13] C. Florensa, Y. Duan, and P. Abbeel. Stochastic neural networks for hierarchical reinforcement
345 learning. *arXiv preprint arXiv:1704.03012*, 2017.
- 346 [14] C. Florensa, D. Held, X. Geng, and P. Abbeel. Automatic goal generation for reinforcement
347 learning agents. In *International Conference on Machine Learning*, pages 1515–1528, 2018.
- 348 [15] K. Gregor, D. J. Rezende, and D. Wierstra. Variational intrinsic control. *arXiv preprint*
349 *arXiv:1611.07507*, 2016.
- 350 [16] S. Gu, E. Holly, T. Lillicrap, and S. Levine. Deep reinforcement learning for robotic manipula-
351 tion with asynchronous off-policy updates. In *2017 IEEE international conference on robotics*
352 *and automation (ICRA)*, pages 3389–3396. IEEE, 2017.
- 353 [17] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy
354 deep reinforcement learning with a stochastic actor. *CoRR*, abs/1801.01290, 2018.
- 355 [18] S. Hansen, W. Dabney, A. Barreto, D. Warde-Farley, T. Van de Wiele, and V. Mnih. Fast task
356 inference with variational intrinsic successor features. In *International Conference on Learning*
357 *Representations*, 2019.
- 358 [19] E. Hazan, S. Kakade, K. Singh, and A. Van Soest. Provably efficient maximum entropy
359 exploration. In *International Conference on Machine Learning*, pages 2681–2691, 2019.
- 360 [20] R. Houthoofd, X. Chen, Y. Duan, J. Schulman, F. De Turck, and P. Abbeel. Vime: variational
361 information maximizing exploration. In *Proceedings of the 30th International Conference on*
362 *Neural Information Processing Systems*, pages 1117–1125, 2016.

- 363 [21] T. Jaksch, R. Ortner, and P. Auer. Near-optimal regret bounds for reinforcement learning.
364 *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.
- 365 [22] Y. Jinnai, J. W. Park, D. Abel, and G. Konidaris. Discovering options for exploration by
366 minimizing cover time. In *International Conference on Machine Learning*, pages 3130–3139.
367 PMLR, 2019.
- 368 [23] Y. Jinnai, J. W. Park, M. C. Machado, and G. Konidaris. Exploration in reinforcement learning
369 with deep covering options. In *International Conference on Learning Representations*, 2020.
- 370 [24] L. Lee, B. Eysenbach, E. Parisotto, E. Xing, S. Levine, and R. Salakhutdinov. Efficient
371 exploration via state marginal matching. *arXiv preprint arXiv:1906.05274*, 2019.
- 372 [25] S. H. Lim and P. Auer. Autonomous exploration for navigating in MDPs. In *Conference on*
373 *Learning Theory*, pages 40–1, 2012.
- 374 [26] Y. Liu and E. Brunskill. When simple exploration is sample efficient: Identifying sufficient
375 conditions for random exploration to yield pac rl algorithms. *arXiv preprint arXiv:1805.09045*,
376 2018.
- 377 [27] K. Lu, A. Grover, P. Abbeel, and I. Mordatch. Reset-free lifelong learning with skill-space
378 planning. *arXiv preprint arXiv:2012.03548*, 2020.
- 379 [28] M. C. Machado, M. G. Bellemare, and M. Bowling. A laplacian framework for option discovery
380 in reinforcement learning. In *International Conference on Machine Learning*, pages 2295–2304.
381 PMLR, 2017.
- 382 [29] M. C. Machado, C. Rosenbaum, X. Guo, M. Liu, G. Tesauro, and M. Campbell. Eigenoption
383 discovery through the deep successor representation. In *International Conference on Learning*
384 *Representations*, 2018.
- 385 [30] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves,
386 M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep rein-
387 forcement learning. *nature*, 518(7540):529–533, 2015.
- 388 [31] N. Modhe, P. Chattopadhyay, M. Sharma, A. Das, D. Parikh, D. Batra, and R. Vedantam.
389 Ir-vic: Unsupervised discovery of sub-goals for transfer in rl. In *Proceedings of the Twenty-*
390 *Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*. International Joint
391 Conferences on Artificial Intelligence Organization, 2020.
- 392 [32] S. Mohamed and D. J. Rezende. Variational information maximisation for intrinsically motivated
393 reinforcement learning. In *Advances in neural information processing systems*, pages 2125–
394 2133, 2015.
- 395 [33] M. Mutti, L. Pratissoli, and M. Restelli. A policy gradient method for task-agnostic exploration.
396 *arXiv preprint arXiv:2007.04640*, 2020.
- 397 [34] V. H. Pong, M. Dalal, S. Lin, A. Nair, S. Bahl, and S. Levine. Skew-fit: State-covering
398 self-supervised reinforcement learning. In *International Conference on Machine Learning*,
399 2020.
- 400 [35] M. L. Puterman. *Markov Decision Processes.: Discrete Stochastic Dynamic Programming*.
401 John Wiley & Sons, 1994.
- 402 [36] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
403 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 404 [37] R. Sekar, O. Rybkin, K. Daniilidis, P. Abbeel, D. Hafner, and D. Pathak. Planning to explore
405 via self-supervised world models. In *International Conference on Machine Learning*, pages
406 8583–8592. PMLR, 2020.
- 407 [38] A. Sharma, S. Gu, S. Levine, V. Kumar, and K. Hausman. Dynamics-aware unsupervised
408 discovery of skills. In *International Conference on Learning Representations*, 2020.

- 409 [39] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker,
410 M. Lai, A. Bolton, et al. Mastering the game of go without human knowledge. *nature*,
411 550(7676):354–359, 2017.
- 412 [40] J. Tarbouriech, E. Garcelon, M. Valko, M. Pirotta, and A. Lazaric. No-regret exploration in
413 goal-oriented reinforcement learning. In *International Conference on Machine Learning*, 2020.
- 414 [41] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012*
415 *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033,
416 2012.
- 417 [42] D. Warde-Farley, T. Van de Wiele, T. Kulkarni, C. Ionescu, S. Hansen, and V. Mnih. Unsuper-
418 vised control through non-parametric discriminative rewards. In *International Conference on*
419 *Learning Representations*, 2019.
- 420 [43] K. Xie, H. Bharadhwaj, D. Hafner, A. Garg, and F. Shkurti. Skill transfer via partially amortized
421 hierarchical planning. In *International Conference on Learning Representations*, 2021.
- 422 [44] K. Xu, S. Verma, C. Finn, and S. Levine. Continual learning of control primitives: Skill
423 discovery via reset-games. *arXiv preprint arXiv:2011.05286*, 2020.
- 424 [45] D. Yarats, R. Fergus, A. Lazaric, and L. Pinto. Reinforcement learning with prototypical
425 representations. *arXiv preprint arXiv:2102.11271*, 2021.
- 426 [46] J. Zhang, H. Yu, and W. Xu. Hierarchical reinforcement learning by discovering intrinsic
427 options. In *International Conference on Learning Representations*, 2021.

428 **Checklist**

- 429 1. For all authors...
- 430 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
431 contributions and scope? [Yes]
- 432 (b) Did you describe the limitations of your work? [Yes] We discuss the limitations and
433 directions of further investigation in the conclusion.
- 434 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We do
435 not foresee any obvious negative societal impacts from our work, which focuses on the
436 fundamentals of reinforcement learning and proposes a new algorithm for unsupervised
437 skill discovery.
- 438 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
439 them? [Yes]
- 440 2. If you are including theoretical results...
- 441 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 442 (b) Did you include complete proofs of all theoretical results? [N/A]
- 443 3. If you ran experiments...
- 444 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
445 imental results (either in the supplemental material or as a URL)? [No] We plan to
446 release our code upon acceptance of this work.
- 447 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
448 were chosen)? [Yes] See Appendix.
- 449 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
450 ments multiple times)? [Yes] Yes when possible.
- 451 (d) Did you include the total amount of compute and the type of resources used (e.g., type
452 of GPUs, internal cluster, or cloud provider)? [No]
- 453 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 454 (a) If your work uses existing assets, did you cite the creators? [Yes] See Sect. 5.
- 455 (b) Did you mention the license of the assets? [N/A]
- 456 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 457
- 458 (d) Did you discuss whether and how consent was obtained from people whose data you’re
459 using/curating? [N/A]
- 460 (e) Did you discuss whether the data you are using/curating contains personally identifiable
461 information or offensive content? [N/A]
- 462 5. If you used crowdsourcing or conducted research with human subjects...
- 463 (a) Did you include the full text of instructions given to participants and screenshots, if
464 applicable? [N/A]
- 465 (b) Did you describe any potential participant risks, with links to Institutional Review
466 Board (IRB) approvals, if applicable? [N/A]
- 467 (c) Did you include the estimated hourly wage paid to participants and the total amount
468 spent on participant compensation? [N/A]

470 Appendix

472 Table of Contents

473	A UPSIDE Algorithm	14
474	B Environment Details	16
475	C Experimental Details	16
476	D Additional Experiments	19
477	E An interpretation of our optimization problem	22

481 A UPSIDE Algorithm

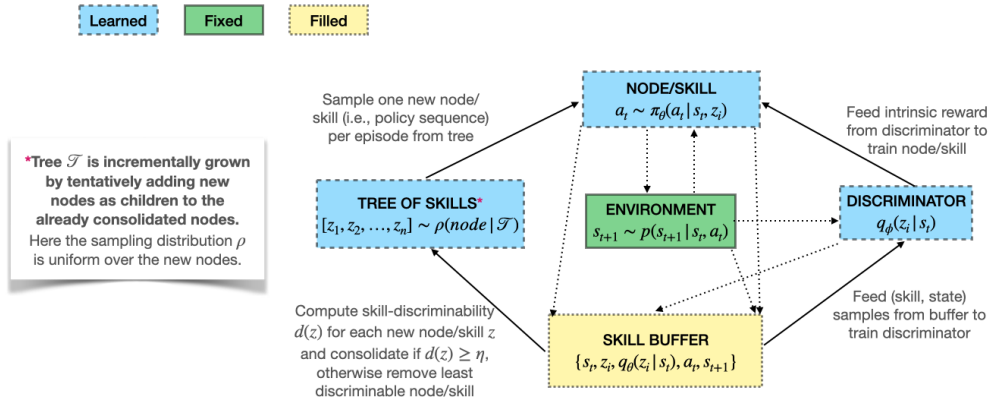


Figure 7: High-level approach of UPSIDE

482 We provide a diagram of the high-level approach of UPSIDE in Fig. 7 and a detailed pseudo-code in
 483 Alg. 2. UPSIDE initializes a tree structure \mathcal{T} with root node 0 and queue of parent nodes $\mathcal{W} = \{0\}$.
 484 As long as the queue is not empty, the following steps are performed:

- 485 • **(Generating new skills)** We expand the tree at a leaf $w \in \mathcal{W}$ by adding N_0 new nodes/skills
 486 denoted by $\mathcal{C}(w)$ (lines 1, 2). We initialize a discriminator with $|\mathcal{T}|$ classes to account for the newly
 487 created nodes (line 3).
- 488 • **(Skill Learning)** Then the new skills and discriminator are optimized as follows:
 - 489 - We sample (uniformly) and rollout the new skills $z \in \mathcal{C}(w)$ and add the states of their diffusing
 490 parts in their corresponding buffers \mathcal{B}_z (lines 8 to 11).
 - 491 - We update the discriminator by leveraging the states and skill labels in the buffers (lines 12 to
 492 17). In particular, the ratio μ signifies that we train more often the discriminator on the previously
 493 consolidated skills/classes than on the new skills/classes in $\mathcal{C}(w)$, i.e., we sample pairs (label z ,
 494 state in \mathcal{B}_z) with probability $(\mathbb{I}[z \in \mathcal{C}(w)] + \mu \mathbb{I}[z \notin \mathcal{C}(w)]) / ((1 - \mu)|\mathcal{C}(w)| + \mu|\mathcal{T}|)$. We give
 495 slightly more weight to the already consolidated skills in the discriminator training because the

Algorithm 2: UPSIDE

Initialize: Discriminability threshold $\eta \in (0, 1)$, branching factor $N_0 \geq 1$ (to be adapted at each node), (optional) maximum branching factor $N_{\max} \geq N_0$, patience K , window \mathcal{V} , number of trajectory rollouts R per update of discriminator and policies, batch size of N_{discr} to train the discriminator, ratio μ of probabilities between consolidated classes and new classes to train discriminator

Initialize: Tree \mathcal{T} initialized as a root node indexed by 0, queue of parent nodes $\mathcal{W} = \{0\}$.

```
while  $\mathcal{W} \neq \emptyset$  do // tree expansion
1  Dequeue a node/skill  $w \in \mathcal{W}$  and expand  $\mathcal{T}$  at  $w$  by adding a set  $\mathcal{C}(w)$  of  $N_0$  nodes/skills
2  Create random policies  $\pi_z$  and buffers  $\mathcal{B}_z, \forall z \in \mathcal{C}(w)$ 
3  Initialize discriminator  $q_\phi$  with  $|\mathcal{T}|$  classes
4  Continue = true; Saturated = false
5  while Continue do
6      for  $K$  iterations do
7          for  $r \in [1, R]$  do // collect  $R$  trajectories
8              Sample a skill  $z$  from  $\mathcal{C}(w)$  at random
9              Extract the sequence of nodes  $z_{(1)}, \dots, z$  in  $\mathcal{T}$  leading to  $z$ 
10             Execute the composed (directed part) policy  $(\pi_{z_{(1)}}, \dots, \pi_z)$  followed by the diffusing part
11             Add states observed during the diffusing part to  $\mathcal{B}_z$ 
12              $B = \{\}$  // Initialize batch to update the discriminator
13             while  $|B| < N_{\text{discr}}$  do
14                 Sample a skill  $z$  from  $\mathcal{T}$  w.p.  $(\mathbb{I}[z \in \mathcal{C}(w)] + \mu \mathbb{I}[z \notin \mathcal{C}(w)]) / ((1 - \mu)|\mathcal{C}(w)| + \mu|\mathcal{T}|)$ 
15                 Sample a state  $s$  from the last  $\mathcal{V}$  states of  $\mathcal{B}_z$ 
16                 Add  $(s, z)$  to  $B$ 
17             Update discriminator  $q_\phi$  with SGD on  $B$  to predict label  $z$ 
18             for  $z \in \mathcal{C}(w)$  do
19                 Update policy  $\pi_z$  using SAC to optimize the discriminator reward as in Sect. 3.1.
20             Compute the skill-discriminability  $d(z) = \hat{q}_\phi^{(B)}(z) = \frac{1}{|B|} \sum_{(s,z) \in B} q_\phi(z|s)$  for all  $z \in \mathcal{C}(w)$ 
21             if  $\min_{z \in \mathcal{C}(w)} d(z) < \eta$  then // Node removal
22                 Remove the node/skill  $z = \arg \min_{z \in \mathcal{C}(w)} d(z)$  from  $\mathcal{C}(w)$  and  $\mathcal{T}$ 
23                 Set Saturated = true
24             else if not Saturated then
25                 Add one new node/skill to  $\mathcal{C}(w)$  and  $\mathcal{T}$ 
26                 if  $|\mathcal{C}(w)| = N_{\max}$  then
27                     Set Saturated = true
28             else
29                 Set Continue = false
30             Enqueue in  $\mathcal{W}$  the consolidated nodes  $\mathcal{C}(w)$ 
```

496 discriminator is reinitialized whenever new classes (i.e., nodes) are added, thus we seek to avoid the
497 new classes from invading the territory of the older classes that were previously correctly learned.
498 In addition, we only update the discriminator on recent batches of data from the buffers via the
499 window \mathcal{V} (which considers only the last \mathcal{V} states in each skill buffer), which is more sample
500 efficient than doing the discriminator update in a fully on-policy manner (e.g., [12]), especially
501 in our setting where the discriminator changes over training as new skill-nodes (i.e., classes) are
502 added.

503 - We update the policies of the new skills/nodes in $\mathcal{C}(w)$ with SAC to optimize the intrinsic
504 reward of the discriminator predictions as explained in Sect. 3.1 (line 19). Note that we keep fixed
505 the policies of the previously consolidated nodes/skills, which makes the learning of the tree more
506 stable.

507 • **(Node Consolidation)** After a *patience* period characterized by K iterations of training (line 6),
508 if all skills are η -consolidated (i.e., the constraint of problem (5) is verified), we tentatively add
509 more skills to the leaf w (line 25). On the other hand, if any skill does not meet the discriminability
510 threshold, we remove it and seek to consolidate the remaining skills into the tree (line 22). The role
511 of the Saturated and Continue booleans is to ensure that the node addition operation cannot be
512 performed if a node removal operation has already been performed in the training of the set $\mathcal{C}(w)$.
513 Recall that the function is monotone, so if a skill is removed, the optimum cannot be larger. The
514 (optional) N_{\max} value represents the maximum branching factor (i.e., number of children nodes)
515 imposed at each node of the tree.

516 B Environment Details

517 **Continuous mazes.** We consider mazes with height and width 50. The state space is continuous,
518 and there are some horizontal and verticals walls of width 1. The agent observes its current (x, y)
519 Cartesian position (i.e., it does not observe the walls) and it outputs actions $[dx, dy]$ that control its
520 location. The actions dx and dy are constrained to be in $[-1, +1]$. The movement of the agent is
521 affected by collisions with walls: when the agent collides with a wall, it stays in its original position.

522 **CartPole.** We modify slightly the simulator from OpenAI Gym [8] to make exploration more
523 difficult and thus to make it more challenging to learn diverse behaviors: the agent moves along the x
524 horizontal position between -2.4 and 2.4 and the pole starts in the reverse position at $x = 0$. When
525 the agent goes out of the x interval, it is teleported back to its initial position (but there is no reset).
526 Observations are $(x, \dot{x}, \theta, \dot{\theta})$ where x is the horizontal position, θ is the angle of the pole to the x -axis,
527 and $\dot{x}, \dot{\theta}$ are their respective velocities.

528 **Reacher.** We use the standard MuJoCo implementation of Reacher [41], which is a two-joint
529 robotic arm where the action space $([-1, +1])$ is the torque applied to both joints with gear 30.

530 C Experimental Details

531 C.1 Baselines

532 For all methods, we augment the state space with the current time-step because horizons are finite.

533 **DIAYN-K.** This corresponds to the original DIAYN algorithm [12] where K is the number of skills
534 to be learned. In order to make the architecture more similar to UPSIDE, we use distinct policies for
535 each skill, i.e. they do not share weights as opposed to [12]. While this may come at the price of
536 sample efficiency, it may also help put lesser constraint on the model (e.g. gradient interference).

537 **DIAYN-Curriculum.** We augment DIAYN with a curriculum that enables to be less dependant
538 on an adequate tuning of the hyperparameter of the number of skills of DIAYN. We consider the
539 curriculum of UPSIDE where we start from either a large or small number N_0 of skills, learn skills
540 during a period of time/number of interactions. If the configuration satisfies the discriminability
541 threshold η , a skill is added, otherwise a skill is removed or learning stopped (as in Alg. A, lines
542 20-29). Note that the increasing version of this curriculum is similar to the one proposed in VALOR [1,
543 Sect. 3.3].

544 **SMM.** We used SMM [24] as it is state-of-art in terms of coverage, at least on long-horizon control
545 problems, although [10] reported poor performance in hard-to-explore bottleneck mazes. We tested
546 the regular SMM version, i.e. learning a state density model with a VAE, yet we failed to make it work
547 on the Mazes domain. As we use the cartesian xy positions in maze domains, learning the identity
548 function on two-dimensional input data is too easy with a VAE, thus preventing the benefits of using
549 a density model to drive exploration. Thus we considered a more straightforward implementation of
550 SMM by using the “real” state distribution through counting. Specifically, we maintain a discretized
551 state distribution by counting states in buckets (similar to the way we compute the achieved coverage).
552 The distribution is just computed by dividing by the sum over buckets. We did not use a moving
553 average so counts are not forgotten: the state distribution is over all policies encountered since the
554 beginning of training (whereas the state distribution is “online” in [24]).

555 C.2 Architecture and Hyperparameters

556 The architecture of the different methods remains the same in all our experiments, except that the
557 number of hidden units changes across considered environments. We consider decoupled actor and
558 critic in SAC, they both have the same (but unshared weights) state processing architectures. The
559 observation and the step are passed through non-linear MLP with 1 hidden layer with units h , then
560 are concatenated. The concatenation is then mapped to an embedding. For the actor, this embedding
561 is mapped to a mean and variance embedding, then passes through a Squashed Gaussian as explained

562 in [17]. For the critic, the embedding is concatenated with a non-linear (1 hidden layer) embedding
563 of the action, then passed through a final non-linear MLP (1 layer) to a one-dimensional value.

564 The discriminator is a two-hidden layer model with output size the number of skills in the tree.

565 **Environment-specific hyperparameters.** Mazes: $h = \{16, 64\}$ hidden units per layer for policy,
566 and $h = 128$ hidden units per layer for discriminator. Continuous control domains: $h = 256$ hidden
567 units per layer for both policy and discriminator.

568 **Common (for methods and environments) optimization hyperparameters.** (See App. A for
569 meaning of each hyperparameter)

- 570 • SAC entropy: $\{0.1, 0.01, 0.001\}$
- 571 • discount factor: $\gamma = 0.99$
- 572 • Q-function soft updates $\tau = 0.005$
- 573 • learning rates $lr_{\text{policy}} = 0.001$, $lr_{\text{discriminator}} = \{0.0001, 0.001\}$
- 574 • discriminator batch size $B = 1024$
- 575 • $\mu = \{2, 5\}$
- 576 • $\mathcal{V} = 100$
- 577 • Replay buffer size: $1e6$

578 Note that hyperparameters are kept fixed for the downstream tasks too.

579 For UPSIDE and DIAYN-curriculum, we set the patience to be a time-limit instead of a number of
580 iterations. We tried both 300 and 600 seconds to avoid the running time getting too high if the tree
581 grows large.

582 The total running time for DIAYN-K and SMM is the same than the maximum running time of UPSIDE.

583 C.3 Model selection

584 We train all methods with a grid search over the set of hyperparameters described in App. C.2, for
585 multiple seeds, which we call *unsupervised seeds*, to evaluate robustness over both the initialization
586 of model weights and randomness of the algorithm. For each unsupervised seed, we select the
587 set of hyperparameters that has maximum value for the criterion of number of skills for UPSIDE,
588 DIAYN-curriculum and for the criterion of average discriminability for DIAYN-K and SMM.

589 With this set of hyperparameters per seeds, we can then report some measurement, e.g. coverage,
590 averaged over unsupervised seeds.

591 C.4 Evaluation protocol

- 592 1. We train the method in its unsupervised phase.
- 593 2. We then do model selection as explained in App. C.3, which gives a model per method per
594 unsupervised seed.
- 595 3. We rollout N episodes per model and compute coverage as explained in the main paper in Sect. 5.
596 Coverage is averaged over unsupervised seeds.
- 597 4. For each model (associated to a method) and unsupervised seed, we run the downstream tasks (as
598 explained in App. C.5), with the same grid search over hyperparameters, with additional seeds,
599 which we call *downstream seeds*.
- 600 5. For each method and unsupervised seed, we do model selection over downstream seeds on the
601 criterion of reward.
- 602 6. We plot the reward averaged over unsupervised and downstream seeds, with error bars for each
603 method.

C.5 Downstream task scenario in Mazes

We consider the downstream task of quickly finding and then reliably reaching an unknown goal, summarized in Alg. 3. There exists a goal region \mathcal{G} with unknown coordinates $(x_{\mathcal{G}}, y_{\mathcal{G}})$ that can be identified only once it is reached. The unknown nature of the goal and its sparse identification signal (i.e., reward $r_{\mathcal{G}}(s) = \mathbb{1}[s \in \mathcal{G}]$) makes the problem challenging, as the agent must perform “blind” and exhaustive exploration so as to encounter the goal as quickly as possible. UPSIDE’s clustering of the state space with its ability to navigate efficiently to any given cluster is a desirable property to tackle this problem. In Alg. 3, we uniformly sample the nodes of the tree (i.e., execute the diffusing part of each skill) until the goal is found. Note that we use a budget of K iterations (which could be either environment interactions or time) for UPSIDE to find the goal with the tree, otherwise we train a policy with SAC on the reward.

604

Once the goal is identified, this becomes a standard goal-oriented task, where no distance-to-goal is available, i.e., the reward signal is *sparse*, which makes the learning problem more difficult. The design of UPSIDE enables to identify the closest skill to the goal according to the learned discriminator, and we then fine-tune its diffusing part into a goal-oriented policy, as shown in Alg. 4.

The same approach is used for DIAYN and SMM. For SAC, a plain policy is trained directly on the reward signal.

We thus see that this task calls for a dual property of coverage and directedness.

Goals g were sampled uniformly in the available state-space, but for the sake of simplicity, we only show in Section 5 two representative goal positions, a moderately close goal and a far goal. The goal region is a circle with radius 1, thus the agent gets rewarded 1, when $\|s - g\|_2^2 < 1$.

Algorithm 3: Unknown goal

Input: Unknown goal region \mathcal{G} ,
Budget K .

```

for  $K$  iterations do // Find  $\mathcal{G}$ 
  Sample  $z$  in  $\mathcal{T}$  at random
  Extract the sequence of nodes
   $z_{(1)}, \dots, z$  in  $\mathcal{T}$  leading to  $z$ 
  Execute the composed
  (directed part) policy
   $(\pi_{z_{(1)}}, \dots, \pi_z)$  followed by
  the diffusing part of  $z$ 
  Stop if  $\mathcal{G}$  is reached.
if  $\mathcal{G}$  was found then
  | Run Alg. 4 with goal  $\mathcal{G}$ 
else
  | Train SAC policy on the
  | reward

```

Algorithm 4: Known goal

Input: Known goal region \mathcal{G} .
Compute skill-node

$$z^* = \arg \max_{z \in \mathcal{T}} \sum_{g \in \mathcal{G}} q_{\phi}(z|g).$$

Fine-tune the diffusing part of
skill-node z^* via RL with reward
 $r_{\mathcal{G}}(s) = \mathbb{1}[s \in \mathcal{G}]$.

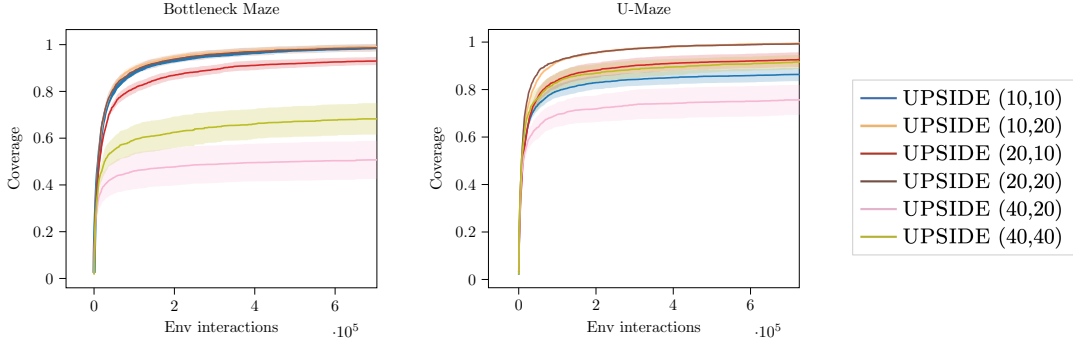


Figure 8: Ablation on the values of T and H for UPSIDE on the bottleneck and U-maze.



Figure 9: Difficulty of UPSIDE to cover the mazes if the hyperparameters T, H are set too large w.r.t. the environment size (here, $T = H = 40$, and we recall that the mazes are of size 50×50). Top (resp. bottom) row corresponds to the stochastic (resp. deterministic) executions of the policies of the directed parts of the skills.

605 D Additional Experiments

606 In this section, we report additional experiments. We ran all methods with 3 unsupervised seeds for
 607 each set of hyperparameters. All plots are generated according to the evaluation protocol explained in
 608 App. C.4.

609 D.1 Ablation on the skill lengths T and H

610 We investigate the sensitiveness of UPSIDE w.r.t. T and H , the lengths of the directed and diffusing
 611 part of the skill, respectively. Fig. 8 shows that the method is quite robust to reasonable choices
 612 of T and H , although there exists configurations where UPSIDE does not achieve full coverage, in
 613 particular in the bottleneck maze when T and H are too large (e.g., $T = 40, H = 20$), see also
 614 Fig. 9. This makes sense as the environments require “narrow” exploration (e.g., the bottleneck region
 615 that the agent must “escape” from is quite small), thus composing disproportionately long skills
 616 may hinder the coverage. Moreover, increasing T and H makes the RL training longer and more
 617 challenging (e.g., the reward is more delayed).

618 D.2 Visual example how the tree learned by UPSIDE fits the environment

619 We investigate the adaptivity (w.r.t. the input branching factor) of the tree structure of UPSIDE
 620 and illustrate that it can properly fit the unknown environment. As demonstrated in 10, UPSIDE
 621 successfully covers a large part of the tree maze, which is quite hard to explore given its narrow
 622 corridors. Here $T = 5$ and $H = 10$, and the branching factor N_0 is set to 3. In the terminal region



Figure 10: (Left) Unbalanced tree-shaped maze and (Right) the tree structure learned by UPSIDE. We see that it can successfully *map* the underlying structure of the unknown environment.

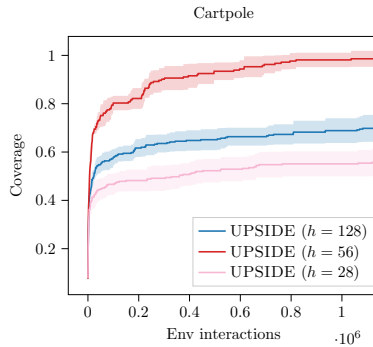


Figure 11: Ablation on h , the number of hidden units per layer of the discriminator of UPSIDE

623 of skill 1 (yellow), it is crucial to consolidate two skills 2 and 3 so that the tree can grow in both
 624 directions. While the tree may have expanded two skills 4 and 5 straight from 3, we see that the skill
 625 4 (blue) overlaps with the intersection of the two small corridors, thus it is the only one sufficiently
 626 discriminable at this tree level, and UPSIDE covers the bottom right corridor in the subsequent level
 627 (i.e., from skill 4 to skill 5 in purple).

628 D.3 Ablation on the capacity of the discriminator

629 On CartPole, we found that it was quite easy for the discriminator to separate skills, though they had
 630 close behaviors “visually”. This can be explained by the fact that high-dimensional states are easier
 631 to discriminate. By reducing the capacity of the discriminator, skills would be naturally forced to be
 632 more “diverse” and avoid overfitting to certain state space regions. To verify this claim, we perform
 633 an ablation on the number of hidden units per layer of the discriminator (Fig. 11), which reveals that
 634 there is a sweet spot of hidden size where coverage is the best. When the hidden size h is too big (128
 635 or 256 in the main paper), many skills (more) are consolidated, but not diverse in their behavior, thus
 636 the coverage is not that large. On the other hand, when h is too small, it is too hard to discriminate
 637 between skills.

638 D.4 Results with more unsupervised seeds

639 In Fig. 12 we add results with 3 new unsupervised seeds per method and set of hyper-parameters. This
 640 complements Fig. 4 and 5 from the main paper by adding error bars. Compared to the main paper,
 641 training time was increased, thus explaining the slight differences in performance (e.g., UPSIDE,
 642 DIAYN-50, SMM-50 improve compared to the random policy thanks to training time).

643 D.5 Average discriminator performance on Mazes and Cartpole

644 Fig. 13 reports the average discriminability of the skills (UPSIDE, DIAYN-curriculum and DIAYN-50)
 645 during training in Bottleneck maze, U-maze and Cartpole. We make the same observation as for
 646 Reacher (see rightmost plot of Fig 5). The DIAYN-50 skills (green) suffer from low discriminability,

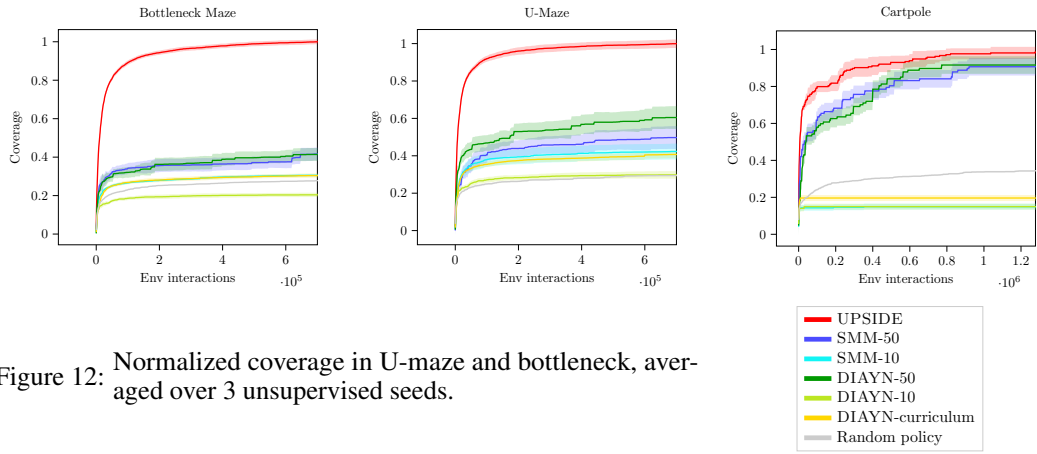


Figure 12: Normalized coverage in U-maze and bottleneck, averaged over 3 unsupervised seeds.

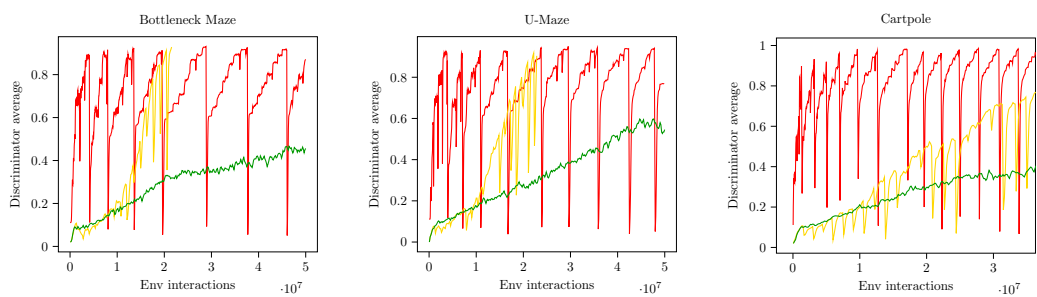


Figure 13: Average discriminability of the skills during training in Bottleneck maze, U-maze and Cartpole

647 while UPSIDE (red) (as well as DIAYN-curriculum in yellow) achieves much higher discriminability
 648 by removing redundant skills.

649 E An interpretation of our optimization problem

650 In this section we provide a theoretically grounded interpretation of the optimization problem solved
 651 by UPSIDE in Sect. 3 and the extent to which it allows to tackle two downstream scenarios: known goal
 652 (Alg. 4) or unknown goal (Alg. 3). Throughout App. E, we consider that the MDP M is finite-state,
 653 finite-action and communicating [35] (i.e., for every pair of states (s, s') , there exists a deterministic
 654 stationary policy under which s' is accessible from s in finite time with non-zero probability).

655 E.1 Preliminaries

656 First we review some concepts and define notation. For any stationary policy π and pair of states
 657 (s, s') , we denote by $\tau_\pi(s, s')$ the (possibly infinite) random variable of the hitting time of state s'
 658 starting from state s following policy π . We then define the *distance* d_π as the expected hitting time
 659 of policy π , i.e.,

$$\tau_\pi(s, s') := \inf\{t \geq 0 : s_{t+1} = s' \mid s_1 = s, \pi\}, \quad d_\pi(s, s') := \mathbb{E}[\tau_\pi(s, s')],$$

660 where the expectation is w.r.t. the random sequence of states generated by executing π starting from
 661 state s . In addition, given a starting state s_0 and distance d , we define the *Max-Distance* as

$$D(d, M, s_0) := \max_{s \in \mathcal{S}} d(s_0, s).$$

662 We can instantiate two commonly considered distances.

- 663 • First, the *random-walk distance* is $d_{\text{RW}}(s, s') := d_{\pi_{\text{RW}}}(s, s')$, where π_{RW} denotes a uniformly stochastic
 664 policy (i.e., whose executed actions are uniformly distributed, independently of the state of the
 665 MDP). Note that for $d \leftarrow d_{\text{RW}}$, D corresponds to the *cover time* starting from s_0 . This notion of
 666 complexity measures how hard it is, in expectation, to cover the entire state space of the MDP
 667 following a uniformly stochastic policy starting from s_0 . It was studied in e.g., [26, 22, 11],
 668 leveraging graph theory [9].
- 669 • Second, the *shortest-path distance* is $d_{\text{SP}}(s, s') := \min_\pi d_\pi(s, s')$ (the minimum can be taken
 670 over the set of stationary deterministic policies [7]). Note that for $d \leftarrow d_{\text{SP}}$, D corresponds to
 671 the *diameter* of the MDP [21, 40] and characterizes the complexity of navigating the state space
 672 starting from s_0 following the set of shortest-path policies. Note that for any state s' , $d_{\text{SP}}(\cdot, s')$ is
 673 the optimal value function under (undiscounted) reward function $r(s) = -\mathbb{1}[s \neq s']$. As such,
 674 numerous methods in goal-conditioned RL (explicitly or implicitly) target the set of policies that
 675 minimize d_{SP} (or variants of it, such as discounted or horizon-truncated) [e.g., 2, 34].

676 E.2 Interpretation

677 An interpretation of our approach is that it performs **clustering** over the state space based on **two**
 678 **different distance functions**:

- 679 • **A tight and difficult-to-deploy distance** d_\star . The tightest metric possible is to consider $d_\star \leftarrow d_{\text{SP}}$,
 680 the shortest-path distance.
- 681 • **A coarse and cheap-to-deploy distance** d_+ . For example we can consider $d_+ \leftarrow d_{\text{RW}}$, the random-
 682 walk distance.

683 Specifically, our approach can be interpreted as seeking to minimize the intra-cluster distance d_+
 684 while maximizing the inter-cluster distance d_\star . Although d_+ is coarser than d_\star , it had the advantage
 685 of being easier to execute the policy to which it corresponds (e.g., a random policy). The implicit
 686 assumption that we make is that d_+ is a *decent enough proxy* for d_\star for *small horizon*, although it
 687 degrades sharply as the horizon increases.

688 In App. E.3 we analyze a simplified structure of our algorithm UPSIDE which allows us to theoretically
 689 analyze the extent to which UPSIDE can tackle the two downstream scenarios explained in App. C.5,
 690 depending on the environment’s properties. For simplicity we will consider the “flat case” of UPSIDE
 691 (see Rmk. 2 for a discussion on the extension to the tree case). Before analyzing the downstream
 692 scenarios (App. E.4.1), we begin by analyzing the properties of the directed part (App. E.3.1) and the
 693 diffusing part (App. E.4) of each UPSIDE skill.

694 **E.3 An analysis of two downstream scenarios tackled by UPSIDE**

695 **E.3.1 Directed part of each UPSIDE skill**

696 **Structure/Assumptions.**

- 697 • The directed part of each skill k is characterized by a pair (π_k, c_k) where the policy π_k is of length
- 698 T and aims to attain a goal state $c_k \in \mathcal{S}$ (chosen by the skill).
- 699 • From the optimization problem of UPSIDE, each skill is η -consolidated according to the discrim-
- 700 inator. We consider that the latter discriminates between the goals $\{c_k\}_{k \in [K]}$ given the current
- 701 state. We then have that $\min_{k \in [K]} q_\phi(c_k | s_T) \geq \eta$.
- 702 • Finally, we assume that the predictions of the discriminator can serve as $\varepsilon_{\text{discr}}$ -accurate approx-
- 703 imations of the probability of π_k reaching the centroid c_k within its length of T steps, where
- 704 $0 \leq \varepsilon_{\text{discr}} < \eta$. This implicitly assumes that we can connect the discriminability property and the
- 705 directedness property (respectively appearing in the reverse and forward forms of MI).

706 We first notice that the directed part π_k has an intrinsic reward signal that approximately targets a
 707 goal-directed behavior. Indeed, as argued in [34, App. E], having an intrinsic reward signal of $r_z(s)$
 708 scaling as $p(c_k | s)$ would amount to learning a goal-oriented policy with goal c_k . In particular, the
 709 optimal non-episodic policy π^\dagger that minimizes $\mathbb{E}[\sum_{t=1}^{+\infty} (1 + \beta \mathcal{H}(\pi(\cdot | s_t))) \mathbb{1}[s_t \neq c_k]]$ induces a
 710 distance-to-goal of $d_{\pi^\dagger}(s, c_k) \leq (1 + \beta \log A) d_{\text{SP}}(s, c_k)$, i.e., it targets the shortest path up to an
 711 entropy bias. However, algorithmically, the directed parts are of length T and the UPSIDE skills reset
 712 every $T + H$ steps. This episodic nature introduces a bias w.r.t. the optimal shortest-path behavior
 713 that is non-trivial to analyze and bound.

714 We now show that thanks to the constraint in the optimization problem of UPSIDE and by our
 715 assumption on the connection between the discriminability property and the directness property, we
 716 can recover goal-directed properties for each first part of skills output by UPSIDE.

717 **Lemma 1.** *Any pair (π_k, c_k) output by UPSIDE verifies*

$$d_{\pi_k}(s_0, c_k) \leq \frac{T + H + 1 - \eta + \varepsilon_{\text{discr}}}{\eta - \varepsilon_{\text{discr}}}.$$

718 *Proof.* Recall that the skill k is episodic of length $T + H$, i.e., it resets to s_0 every $T + H$ time steps.
 719 We denote by $d_\pi^{(T+H)}$ the total number of steps before reaching either the skill’s centroid or $T + H$
 720 steps, and by $f_\pi^{(T+H)}$ the probability of failure to reach the centroid within $T + H$ steps. Then

$$d_{\pi_k}(s_0, c_k) \stackrel{(i)}{=} \frac{d_{\pi_k}^{(T+H)}(s_0, c_k) + f_{\pi_k}^{(T+H)}(s_0, c_k)}{1 - f_{\pi_k}^{(T+H)}(s_0, c_k)} \stackrel{(ii)}{\leq} \frac{T + H + f_{\pi_k}^{(T)}(s_0, c_k)}{1 - f_{\pi_k}^{(T)}(s_0, c_k)},$$

721 where (i) comes from [25, App. B.3], (ii) uses that $d^{(T+H)} \in [0, T + H]$ and that $f_{\pi_k}^{(T+H)}(s_0, c_k) \leq$
 722 $f_{\pi_k}^{(T)}(s_0, c_k)$. We now approximate the probability of failure of reaching the centroid by using the
 723 predictions of the discriminator (the more expressive the discriminator, the better the approximation):
 724 $|1 - f_{\pi_k}^{(T)}(s_0, c_k) - q_\phi(c_k | s_T)| \leq \varepsilon_{\text{discr}}$. The constraint of our optimization problem ensures that the
 725 pair (π_k, c_k) output by UPSIDE satisfies $q_\phi(c_k | s_T) \geq \eta$. Therefore, it holds that

$$\frac{T + H + f_{\pi_k}^{(T)}(s_0, c_k)}{1 - f_{\pi_k}^{(T)}(s_0, c_k)} \leq \frac{T + H + 1 - q_\phi(c_k | s_T) + \varepsilon_{\text{discr}}}{q_\phi(c_k | s_T) - \varepsilon_{\text{discr}}} \leq \frac{T + H + 1 - \eta + \varepsilon_{\text{discr}}}{\eta - \varepsilon_{\text{discr}}}. \quad (6)$$

726 □

727 Note that given any goal state g , having a policy π with bounded $d_\pi(\cdot, g)$ is non-trivial, since it implies
 728 that it reaches the goal with probability 1 (i.e., that it is proper [7]). Also note that the “worst-case”
 729 discriminability property in the constraint (i.e., $q_\phi(c_k | s_T) \geq \eta$) is crucial to obtain Lem. 1, since it
 730 may not be possible to guarantee it given a discriminability property verified on average (e.g., via a
 731 conditional entropy term in the MI).

732 **E.4 Diffusing part of each UPSIDE skill**

733 **Structure/Assumptions.**

- 734 • The diffusing part of skill k is of length H and is composed of a set of states radiating around
 735 c_k , which acts as a centroid for the cluster of states generated by the diffusing part. Formally, we
 736 consider that there exists $\delta > 0$ such that

$$\text{DIFF}(k) := \{y_k : \mathbb{P}(\tau_+(c_k, y_k) \leq H) \geq \delta\}, \quad (7)$$

737 where $\tau_+(s, s')$ denotes the hitting time following the policy that minimizes $d_+(s, s')$.

- 738 • According to the optimization problem solved by UPSIDE, the *clusters* associated to the K skills
 739 *saturate* the state space, i.e., we cannot consolidate an additional cluster. We propose to write this
 740 condition as

$$\forall s \in \mathcal{S}, \exists k \in [K], \exists y_k \in \text{DIFF}(k), \mathbb{P}(\tau_+(s, y_k) \leq H) \geq \delta, \quad (8)$$

741 otherwise from (7) it would be possible to consolidate an additional cluster with centroid s .

- 742 • Finally, we spell out an assumption on the environment that we make throughout App. E.4:

743 **Assumption 1.** *There exists $\Theta \geq 0$ such that*

$$\forall (s, s'), \quad \mathbb{P}(\tau_+(s, s') \leq H) \geq \delta \implies d_*(s, s') \leq H + \Theta.$$

744 *This formalizes the assumption commonly made in goal-conditioned deep RL — either implicitly or*
 745 *explicitly [e.g., 14, Sect. 3.3] — that if a goal is reachable, then there exists a policy that does so*
 746 *reliably. Note that in the special case of a deterministic MDP we have $\Theta = 0$.*

747 **Definition 1.** *We define the following “local” quantities:*

- 748 • For any $s \in \mathcal{S}$ and any $k \in [K]$, define $\Delta_+(s; k) := \max_{y_k \in \text{DIFF}(k)} |d_+(s, y_k) - d_+(y_k, s)|$.

- 749 • For any $s \in \mathcal{S}$ and any $k \in [K]$, define $\Delta_*(s; k) := \max_{y_k \in \text{DIFF}(k)} |d_*(s, y_k) - d_*(y_k, s)|$.

750 *Note that under the communicating MDP assumption, both quantities are always bounded. They*
 751 *measure the level of “reversibility” of the MDP w.r.t. the d_+ and d_* distance, respectively. Moreover,*
 752 *in the special case of an MDP with locally symmetric actions, the distance d_* is symmetric so $\Delta_* = 0$.*

753 We first derive two lemmas and then position their statements w.r.t. the two downstream objectives of
 754 UPSIDE.

755 **Lemma 2.** *It holds that*

$$\forall s \in \mathcal{S}, \exists k \in [K] : \quad \mathbb{P}\left(\tau_+(c_k, s) \leq H + \frac{H + \Delta_+(s; k) + \Theta}{1 - \delta}\right) \geq \delta^2.$$

756 **Lemma 3.** *It holds that*

$$\forall s \in \mathcal{S}, \exists k \in [K] : \quad d_*(c_k, s) \leq 2(H + \Theta) + \Delta_*(s; k),$$

757 *Proof of Lem. 2.* We prove the result by contradiction. Assume the contrary of Lem. 2; then there
 758 exists a state $s \in \mathcal{S}$ such that for every $k \in [K]$, $\mathbb{P}(\tau_+(c_k, s) \leq H + Z_k) < \delta^2$, with $Z_k :=$
 759 $(H + \Delta_+(s; k) + \Theta)/(1 - \delta)$. We now use that the diffusing part of each skill k radiating around its
 760 centroid c_k is composed of states $\{y_k : \mathbb{P}(\tau_+(c_k, y_k) \leq H) \geq \delta\}$. This means that for every $k \in [K]$
 761 and $y_k \in \text{DIFF}(k)$, $\mathbb{P}(\tau_+(y_k, s) \leq Z_k) < \delta$. Noticing that $d_+(y_k, s) = \mathbb{E}[\tau_+(y_k, s)]$ by definition,
 762 we get $d_+(y_k, s) > (1 - \delta)Z_k \geq H + \Theta + d_+(y_k, s) - d_+(s, y_k)$, where the last inequality comes
 763 from the definition of $\Delta_+(s; k)$. Therefore, $d_+(s, y_k) > H + \Theta$. So by contraposition of Asm. 1,
 764 $\mathbb{P}(\tau_+(s, y_k) \leq H) < \delta$. Since this is true for all $k \in [K]$ and $y_k \in \text{DIFF}(k)$, we get a contradiction
 765 on condition (8). \square

766 *Proof of Lem. 3.* Take any state $s \in \mathcal{S}$. Case 1: $\exists k \in [K], s \in \text{DIFF}(k)$. Then $\mathbb{P}(\tau_+(c_k, s) \leq$
 767 $H) \geq \delta$. From Asm. 1 this means $d_*(c_k, s) \leq H + \Theta$. Case 2: $\forall k \in [K], s \notin \text{DIFF}(k)$. Then
 768 from condition (8), there exists $k \in [K]$ and $y_k \in \text{DIFF}(k)$ such that $\mathbb{P}(\tau_+(s, y_k) \leq H) \geq \delta$,
 769 which implies that $d_*(s, y_k) \leq H + \Theta$ from Asm. 1. By definition of $\Delta_*(s; k)$, it holds that
 770 $d_*(s_k, y) \leq H + \Theta + \Delta_*(s; k)$. Furthermore, y_k verifies $\mathbb{P}(\tau_+(c_k, y_k) \leq H) \geq \delta$, which means
 771 from Asm. 1 that $d_*(c_k, y_k) \leq H + \Theta$. We conclude by the triangle inequality that $d_*(c_k, y) \leq$
 772 $d_*(c_k, s_k) + d_*(s_k, y) \leq 2(H + \Theta) + \Delta_*(s; k)$. \square

773 **E.4.1 Analysis of two downstream scenarios tackled by UPSIDE**

774 We consider the two downstream tasks detailed in App. C.5: ① finding an unknown goal (Alg. 3) and
 775 ② reliably reaching a known goal (Alg. 4).

776 These downstream scenarios require the ability to efficiently traverse from s_0 to any state s of the
 777 MDP. Ideally we would deploy the policy associated to $d_*(s_0, s)$, i.e., the shortest-path policy, yet it
 778 is difficult to compute. On the other extreme, deploying the random-walk strategy is very easy yet
 779 much more inefficient, since $d_*(s_0, s) \ll d_+(s_0, s)$. Our approach targets the following intermediate
 780 approach.

781 First, we upper bound using the triangle inequality

$$\max_{s \in \mathcal{S}} d_*(s_0, s) \leq \max_{s \in \mathcal{S}} \left\{ \min_{k \in [K]} d_*(s_0, c_k) + d_*(c_k, s) \right\}. \quad (9)$$

782 Under a **zero-shot downstream set-up**, the training objective of UPSIDE seeks to control the follow-
 783 ing upper bound of (9)

$$\max_{s \in \mathcal{S}} \left\{ \min_{k \in [K]} d_*(s_0, c_k) + d_+(c_k, s) \right\}. \quad (10)$$

784 Under a **few-shot downstream set-up**, UPSIDE fine-tunes the diffusing part of the skill to reach the
 785 desired goal state. As such, it targets (9).

786 We now distinguish between the two types of downstream scenarios.

787 **① The unknown-goal downstream task.**

788 From Lem. 2, whatever the unknown goal state s , there exists a skill k whose diffusing part starting
 789 from its centroid c_k can reach s with strictly positive probability, as long as it is executed long enough
 790 (with length depending in particular on the *local* quantity $\Delta_+(s; k)$).

791 As such, Lem. 1 and Lem. 2 prescribe the following *algorithmic strategy*: in a round-robin fashion
 792 over $k \in [K]$, execute the directed part of skill k plus its diffusing part for *increasing* lengths (i.e.,
 793 starting from H and then gradually increasing it). The unknown goal should then be discovered at
 794 some point, specifically within

$$O \left(\frac{1}{\delta^2} \left(\frac{T + H + 1 - \eta + \varepsilon_{\text{discr}}}{\eta - \varepsilon_{\text{discr}}} + \frac{H + \Delta_+(s; k) + \Theta}{1 - \delta} \right) \right)$$

795 time steps, by combining Lem. 1 and 2.

796 **② The known-goal downstream task.**

797 From Lem. 3, for any known goal state $s \in \mathcal{S}$, there exists a skill $k \in [K]$ from which learning to
 798 reach the goal s can be facilitated. Indeed, the shortest-path distance from its centroid c_k to the goal s
 799 depends on the *local* quantity $\Delta_+(s; k)$ (as well as H and Θ).

800 As such, Lem. 1 (with (6)) and Lem. 3 prescribe the following *algorithmic strategy*: first reach the
 801 centroid c_k for which $k \in \arg \max q_\phi(c_k | s)$ (i.e., execute the directed part of skill k), and second
 802 learn to reach s from c_k (by fine-tuning the diffusing part of skill k).

803 **Remark 1.** Inspecting the quantities in Def. 1 and in Asm. 1 allows to characterize the complexity
 804 of the environment in tackling the two types of downstream tasks. In particular, we see that the
 805 complexity is reduced in environments that are close to deterministic (i.e., smaller Θ in Asm. 1)
 806 and that exhibit a “balanced / symmetric” behavior, with the least bottlenecks possible (i.e., smaller
 807 quantities in Def. 1). In addition, the size of the state space \mathcal{S} and the diameter of the MDP implicitly
 808 play a role in the value of the number of clusters K required and in the choice of T , which must be
 809 large enough to ensure in Lem. 1 that $\eta > 0$ holds in the discriminator predictions.

810 **Remark 2.** In the tree case of UPSIDE, T does not have to be large enough (as needed in the flat case)
 811 since the state space may be covered by sequentially composing the directed parts of the skills of
 812 length T . The equations from the flat case would look the same as in the flat case, yet two quantities
 813 would be replaced: the probability of success of reaching the centroid of each cluster of skill-node n
 814 would go from η to $\eta^{d(n)}$ where $d(n)$ is the depth of skill-node n , and the length of the skill would
 815 go from $T + H$ to $d(n)T + H$.

816