

# Intrinsically Motivated Discovery of Temporally Abstract World Model Graphs

**Anonymous authors**

Paper under double-blind review

**Keywords:** Hierarchical RL, model-based RL, exploration.

## Summary

We seek to design reinforcement learning agents that build plannable models of the world that are abstract in both state and time. We propose a new algorithm to construct a *skill graph*; nodes in the skill graph represent abstract states and edges represent skill policies. Previous works that learn a skill graph use random sampling from the state-space and nearest-neighbor search—operations that are infeasible in environments with high-dimensional observations (for example, images). Furthermore, previous algorithms attempt to increase the probability of all edges (by repeatedly executing the corresponding skills) so that the resulting graph is robust and reliable everywhere. However, exhaustive coverage is infeasible in large environments, and agents should prioritize practicing skills that are more likely to result in higher reward. We propose a method to build skill graphs that aids exploration, without assuming state-sampling, distance metrics, or demanding exhaustive coverage.

## Contribution(s)

1. We provide an algorithm that learns a graph-based, plannable abstraction of the environment, even when the observations are high-dimensional; for example, images.

**Context:** Prior work learned graph abstractions by assuming a distance metric over the state-space (for example, [Bagaria et al. \(2021b\)](#); [Lee et al. \(2022\)](#)), and hence cannot be easily applied to environments with image-based observations.

2. We use ideas from Intrinsic Motivation to design a graph construction algorithm that serves as a high-level exploration objective without needing to learn an accurate one-step model of the world—the drive to build the skill graph allows the agent to solve five challenging exploration problems, directly from pixels.

**Context:** Most prior work in model-based exploration either operates in non-image based domains (for example, [Sharma et al. \(2020b\)](#)), or require the agent to learn one-step models (for example, [Hafner et al. \(2022\)](#); [Mendonca et al. \(2021\)](#)).

3. We provide a method for converting a goal-conditioned value function into a plannable abstract world model, which allows us to use dynamic programming to determine which option to execute at each state.

**Context:** [Lo et al. \(2024\)](#) also study this problem, but they assume that option subgoals and initiation regions are provided to the agent; furthermore, they use the model for reward-shaping ([Ng et al., 1999](#)) rather than for option selection at decision time.

# Intrinsically Motivated Discovery of Temporally Abstract World Model Graphs

**Anonymous authors**

Paper under double-blind review

## Abstract

1 We seek to design reinforcement learning agents that build plannable models of the  
 2 world that are abstract in both state and time. We propose a new algorithm to con-  
 3 struct a *skill graph*; nodes in the skill graph represent abstract states and edges repre-  
 4 sent skill policies. Previous works that learn a skill graph use random sampling from  
 5 the state-space and nearest-neighbor search: operations that are infeasible in environ-  
 6 ments with high-dimensional observations (for example, images). Furthermore, previ-  
 7 ous algorithms attempt to increase the probability of all edges (by repeatedly executing  
 8 the corresponding skills) so that the resulting graph is robust and reliable everywhere.  
 9 However, exhaustive coverage is infeasible in large environments, and agents should  
 10 prioritize practicing skills that are more likely to result in higher reward. We show  
 11 that our agent can solve 4 challenging image-based exploration problems more rapidly  
 12 than vanilla model-free RL and state-of-the-art novelty-based exploration. Then, we  
 13 show that the resulting abstract model solves a family of tasks in VISUALPINBALL not  
 14 provided during the agent’s exploration phase.

## 15 1 Introduction

16 Hierarchical reinforcement learning (HRL) (Barto & Mahadevan, 2003) augments the agent’s prim-  
 17 itive actions with temporally extended actions called *options* (Sutton et al., 1999). An agent with  
 18 temporally extended options can plan more effectively because the accuracy of its predictions do  
 19 not degrade as rapidly over time (Sutton et al., 1999; 2011). Furthermore, well-designed options  
 20 empower the agent to explore more efficiently (Machado, 2019). But, where do good options come  
 21 from and how can an agent construct them autonomously solely from interaction with the environ-  
 22 ment? This problem of option discovery (McGovern, 2002) has been studied extensively (Precup,  
 23 2001; McGovern, 2002; Ravindran, 2004; Konidaris, 2011; Bacon, 2018; Machado, 2019), but most  
 24 algorithms do not simultaneously tackle the complementary problem of learning state abstractions  
 25 (Konidaris, 2019). An agent that simultaneously learns both can build a simplified version of the  
 26 environment (Konidaris et al., 2018), which can then be exploited by off-the-shelf planners to yield  
 27 solutions to complex, long-horizon tasks (Machado et al., 2023; Sutton et al., 2022; 2024).

28 One algorithm that simultaneously discovers options and their corresponding state abstractions is  
 29 Deep Skill Graphs (DSG; Bagaria et al., 2021b). DSG learns options that form a graph, where  
 30 the nodes are subgoals and the edges are the policies that navigate between them. To construct  
 31 this graph, DSG first samples a state uniformly at random from the state-space and then selects its  
 32 nearest node  $n_e$  in the graph. Then, it uses planning inside the graph to reach  $n_e$ , after which it moves  
 33 towards the randomly sampled state. This simple combination of random sampling and expanding  
 34 the nearest neighbor node is inspired by motion planning (LaValle, 1998) and provably causes the  
 35 graph to expand towards large unexplored regions of the state-space (Lindemann & LaValle, 2004).  
 36 While effective, DSG makes the following assumptions that limit its applicability: (a) the ability to  
 37 sample meaningful states uniformly at random from the state-space (a state-sampler), and (b) the  
 38 ability to efficiently compute the nearest neighbor node in the graph (a distance metric).

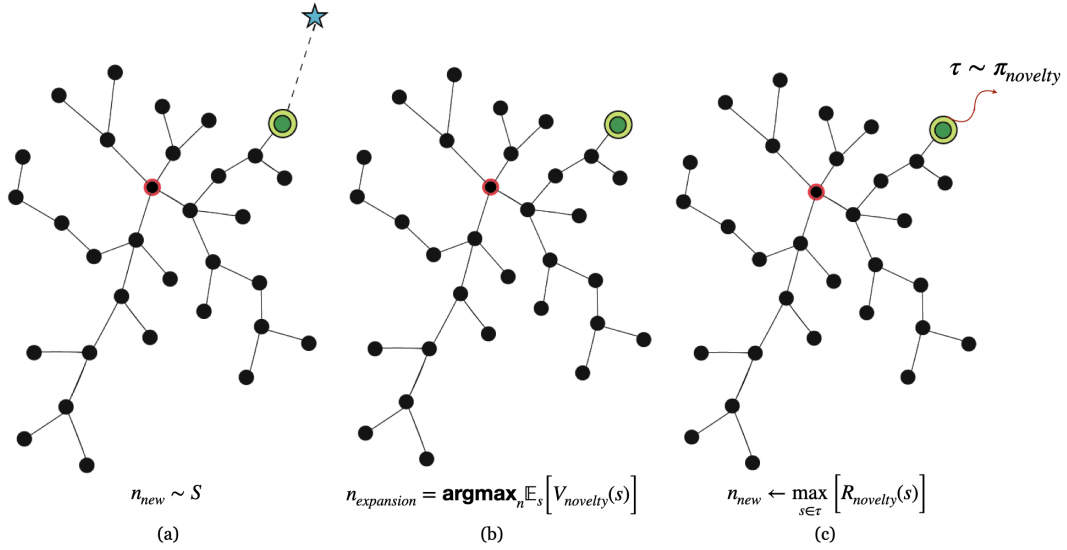


Figure 1: (a) In the original DSG algorithm, a state is sampled uniformly at random (blue star) from the state-space  $S$  and the graph is pulled towards it via its nearest neighbor node (green). (b) In the new algorithm, IM-DSG, the agent uses intrinsic motivation to identify a node to expand using an exploration value function  $V_{novelty}$ . (c) Once the agent reaches the expansion node, it executes an exploration policy  $\pi_{novelty}$ , and the most novel state in the resulting trajectory  $\tau$  is identified as a target for a new skill.

39 How can we incrementally build the skill graph so that it retains its expansion property, but does not  
 40 require a state-sampler or nearest-neighbor search to do so? Instead of generating a random state  
 41 and *pulling* the graph towards it (see Figure 1(a)), we use ideas from intrinsic motivation (Barto,  
 42 2013; Schmidhuber, 2010) to identify an existing node for expansion and *push* the graph outward  
 43 from it (Figure 1(b); Lindemann & LaValle 2004). The agent plans with its existing skills to reach  
 44 the expansion node, and then executes a *low-level exploration policy* trained using RL to maximize  
 45 a novelty-based intrinsic reward (Sutton, 1990; Taïga et al., 2019). Since the intrinsic reward is  
 46 likely to be higher in states that are outside the graph, executing the exploration policy will result  
 47 in a trajectory that will tend to move away from the existing graph, i.e., towards those states where  
 48 the agent does not yet have a reliable abstract model (Figure 1(c)). The agent then identifies the  
 49 most interesting state (Chentanez et al., 2005) in that trajectory and builds a new skill to target it  
 50 (Simsek & Barto, 2004). Over time, the agent refines the skill that reaches this new node in the  
 51 graph, and in doing so, converts this once novel region into an area over which it has gained mastery  
 52 via its abstract model (Veeriah et al., 2018). We call this algorithm Intrinsically Motivated Deep  
 53 Skill Graphs (IM-DSG).

54 We test IM-DSG in a series of exploration problems with image-based observation spaces, to which  
 55 DSG cannot be applied. Our qualitative results show that IM-DSG discovers options that achieve  
 56 intuitively meaningful subgoals in the environment and that the acquired skill graph incrementally  
 57 expands outward from the start state. Quantitatively, we show that our method outperforms pure  
 58 novelty-maximization techniques (Taïga et al., 2019) as well as baseline model-free RL (Kaptur-  
 59 owski et al., 2019) in MINIGRID (Chevalier-Boisvert et al., 2023), VISUAL TAXI (Dietterich, 2000;  
 60 Diuk et al., 2008; Allen et al., 2021), and SOKOBAN (Schrader, 2018); finally we show that the  
 61 learned abstract model is useful in a multi-task context to achieve held-out goals in a zero-shot  
 62 fashion in VISUAL PINBALL (Konidaris & Barto, 2009; Rodriguez-Sanchez & Konidaris, 2024).

## 2 Background and Related Work

We model the interaction of the agent with its environment as a Markov Decision Process  $M = (S, A, R, T, \gamma)$  (Sutton & Barto, 2018). To aid in reward maximization, our agent will discover options  $o \in \mathcal{O} = (\mathcal{I}_o, \pi_o, \beta_o)$  (Sutton et al., 1999), where  $\mathcal{I}_o(s)$  is the probability that  $o$  can be initiated in state  $s$ ,  $\pi_o$  is a policy which outputs primitive actions and  $\beta_o : s \rightarrow \{0, 1\}$  denotes the termination region of the option, which we also call the subgoal of the option. Additionally, each option has a maximum horizon of  $H$  steps and an effect set  $\text{Eff}(o) \subseteq \beta_o$ , which denotes the set of states in which the option actually enters its subgoal  $\beta_o$  (Konidaris et al., 2018).

**Universal Value Function Approximators (UVFAs)** (Schaul et al., 2015) provide a scalable way to learn goal-conditioned value functions  $V_g : S \rightarrow \mathbb{R}$ ,  $Q_g : S \times A \rightarrow \mathbb{R}$  and their corresponding goal-conditioned policies  $\pi_g : S \rightarrow A$  (Kaelbling, 1993) using function approximation. We use UVFAs to parameterize option value functions and policies: instead of representing each option’s policy separately (Sutton et al., 2011), all options can condition the same function approximator on its own subgoals (Bagaria et al., 2021a):  $\pi_o(s) = \text{argmax}_a Q_g(s, a) = \text{argmax}_a Q_{\beta_o}(s, a)$ .

**Exploration in RL.** In tabular RL, uncertainty about the transition and reward models can be captured using visitation counts: if the bonus reward (Sutton, 1990) for action  $a$  from state  $s$  decays as  $1/\sqrt{N[s, a]}$ , the agent can rapidly learn the optimal policy  $\pi^*$  (Lattimore & Szepesvári, 2020; Strehl & Littman, 2008). When the environment is too big to permit tabular counting, counts can be generalized to *pseudocounts*  $N(s)$  (Bellemare et al., 2016): states that are similar to one another have similar pseudocounts. Coin Flip Network (CFN) (Lobel et al., 2023) is a simple, principled, and state-of-the-art technique which uses a neural network  $\phi$  to predict pseudocounts in high-dimensional state-spaces. Novelty seeking agents exhibit impressive performance (Kapturowski et al., 2022), but they do not learn a model suitable for high-level planning. Our algorithm combines novelty maximization (which we refer to as “low-level exploration”) with a drive to learn an abstract, plannable model of the world (“high-level exploration”).

**Option discovery methods.** A large body of work has sought to discover options from the set of primitive actions available to the agent. At a high-level, option discovery techniques can be categorized in the following way:

- **Reward-driven methods.** Option critic (Bacon et al., 2017) and feudal methods (Dayan & Hinton, 1993; Vezhnevets et al., 2017; Nachum et al., 2018b; Levy et al., 2019) learn skills using the overall reward maximizing objective of the agent. Since all learning stems from the extrinsic reward function, they struggle in hard exploration problems with sparse rewards; there are some rare exceptions (McClinton et al., 2021), but they are plagued by the difficulty that the action-space for the high-level policy is the entire state-space of the environment (Nachum et al., 2018a).
- **Empowerment** maximizing methods learn a diverse set of options (Eysenbach et al., 2019b; Zahavy et al., 2021) that increase the agent’s control over its environment (Mohamed & Jimenez Rezende, 2015; Choi et al., 2021; Gregor et al., 2016). These algorithms have shown promise in challenging control (Sharma et al., 2020a) and exploration problems (Baumli et al., 2021; Hansen et al., 2021; Campos Camúñez et al., 2020). Among these, DADS (Sharma et al., 2020a) combines empowerment with high-level planning, but it decomposes the agent’s life into separate training and test phases, i.e., they do not provide a way to *interleave* exploration and model learning. Furthermore, methods like DADS require learning one-step transition models, perhaps explaining why they have not been applied to image-based problems.
- **Spectral methods** use the graph Laplacian (Chung & Graham, 1997) to find the principal directions of the state-space (Mahadevan & Maggioni, 2007) for exploration (Machado et al., 2017). Typically, an option is associated with each eigenvector of the Laplacian; these options demonstrate strong exploratory properties (Jinnai et al., 2019; 2020; Klissarov & Machado, 2023). Both empowerment and spectral methods are promising directions for exploration research, but have not yet been combined with model-based planning, which is our core motivation.

**Deep Model-based RL** methods like Dreamer (Hafner et al., 2018) and MuZero (Schrittwieser et al., 2021) learn a physics-based model of the environment for planning. These algorithms have demonstrated impressive results in a variety of domains (Hafner et al., 2023). However, all learning (even in LEXA (Mendonca et al., 2021) and Director (Hafner et al., 2022), which learn temporal abstractions) is predicated on accurately approximating **one-step** transition and reward models in the agent’s native state-action space. One-step models suffer from the “curse of horizon” (Sutton et al., 1999): model prediction errors compound over time (Talvitie, 2017). Our work seeks to bypass the challenges of model-learning in the raw state-action space: model-free policies are learned in the low-level and temporally-extended models are learned and used in the high-level.

**Go-Explore** shares a similar motivation to our algorithm: exploration is more effective when initiated from the frontiers of the agent’s experience (Ecoffet et al., 2021). But, they require a pre-specified “cell” representation and exploit the near-determinism of Atari domains (Machado et al., 2018) to return to the frontier (Oh et al., 2018). Crucially, Go-Explore does not learn a plannable model while doing exploration, which is the focus of our work.

**Existing graph-based HRL methods** like deep skill graphs (DSG) (Bagaria et al., 2021b; Lee et al., 2022; Lo et al., 2022; 2024) do not generalize to image-based observation spaces. Others do (Konidaris et al. (2018); Lo et al. (2022); Mann et al. (2015); Sutton et al. (2024); Rodriguez-Sanchez & Konidaris (2024), but assume that options are given. Eysenbach et al. (2019a); Huang et al. (2019); Campos Camúñez et al. (2020); Nasiriany et al. (2019); Sharma et al. (2020a) decompose the agent’s training into a pre-training stage and a planning phase: skills are learned over the entire domain and *then* planning is done to improve learned skills. These methods typically perform uniform resets: the agent is randomly reset to a new location at the beginning of every episode (Eysenbach et al., 2019a; Huang et al., 2019; Zhang et al., 2021); this bypasses the exploration problem, which is a key motivation for our work.

### 3 Intrinsically Motivated Deep Skill Graphs

Exploration and planning are deeply intertwined. Exploration aids planning by preferentially seeking out data from regions where the model is currently too weak to support effective planning. But, planning also supports effective exploration: the agent can plan with its abstract model so that it can reach interesting areas reliably, and then conduct local exploration in a more targeted fashion. The Intrinsically Motivated Deep Skill Graphs (IM-DSG) agent interleaves exploration and planning in this way. To support planning, it builds a graph-based model of the environment. Exploration proceeds both inside the graph (to strengthen the existing model) and outside it (to increase the region over which the agent can plan). Inside the graph, the agent refines skills that are currently weak and could lead to high reward; outside the graph, the agent uses a novelty and reward maximizing policy to search for new salient regions that can be reached reliably in the future using its model.

#### 3.1 Semantics of the Skill Graph

Before discussing our algorithm for discovering options and planning with them, we first outline the major components of the skill graph:

**Nodes and Edges.** The skill graph is a weighted and directed graph  $\mathcal{G} = (\mathcal{N}, E, W)$ . Each **node**  $n \in \mathcal{N}$  represents the terminating states of an option  $o_n$ —i.e., its subgoal—determined by the termination condition  $\beta_{o_n} : s \rightarrow \{0, 1\}$ . For each node  $n$ , we can measure the extent to which a state  $s$  belongs to  $n$  by considering the effect of executing option  $o_n$  from any state in the state space:

$$n(s) := \text{Eff}(o_n)(s) = \mathbb{P}(s \mid o_n); \quad (1)$$

that is, the distribution of states after executing  $o_n$ . Moving forward, we slightly abuse notation by taking expectations with respect to the states of a node,  $\mathbb{E}_{s \sim n}[\cdot]$  to mean expectations taken with respect to the effect distribution  $\mathbb{E}_{s \sim \text{Eff}(o_n)}[\cdot]$ .

157 An **edge**  $e_{n_i \rightarrow n_j} \in E$  exists between nodes  $n_i$  and  $n_j$  if executing the option policy  $\pi_{o_{n_j}}$  from  
 158 states in node  $n_i$  is likely to reach node  $n_j$  (Konidaris et al., 2018). The weight of an edge is  
 159 the probability of successfully traversing it in a single option execution, which is equivalent to the  
 160 initiation probability of option  $o_j$  from states inside node  $n_i$ :  $w_{n_i \rightarrow n_j} = \mathbb{E}_{s \sim n_i} [\mathcal{I}_{o_{n_j}}(s)]$  (Bagaria  
 161 et al., 2023).

162 **Mapping from states to nodes and back.** Given a state  $s$ ,  $\mathcal{N}(s)$  denotes the set of nodes that  $s$   
 163 maps to:

$$\mathcal{N}(s) = \{n \in \mathcal{N} \mid \beta_{o_n}(s) = 1\}. \quad (2)$$

164 Of course, it is possible that a state does not satisfy *any* option’s termination condition; so, we aug-  
 165 ment the nodes in the graph with a null vertex  $\varphi$  such that  $\mathcal{N}(s) = \{\varphi\}$  for such states; informally,  
 166  $\varphi$  represents “falling off the graph”.

167 The descendants  $\mathcal{D}(n)$  of a node  $n$  is the set of nodes that can be reached from  $n$  with nonzero  
 168 probability, while staying inside the graph:

$$\mathcal{D}(n) = \{n' \in \mathcal{N} \mid \mathcal{G}.\text{has-path}(n, n')\}, \quad (3)$$

169 where  $\mathcal{G}.\text{has-path}(n, n')$  is a sub-routine implemented using breadth-first search on the adjacency  
 170 matrix representation of the graph  $\mathcal{G}$ . By combining Equations 2 and 3, we can enumerate the nodes  
 171 reachable from the current state  $s_t$  as  $\mathcal{D}(s_t) = \cup_{n \in \mathcal{N}(s_t)} \mathcal{D}(n)$ .

## 172 3.2 Graph Construction Algorithm

173 At a high-level, the graph construction algorithm proceeds as follows: the IM-DSG agent first iden-  
 174 tifies which node in the graph to expand (Section 3.2.1), then it constructs (Section 3.2.4) and solves  
 175 an abstract MDP (Section 3.2.5), which yields a policy over options that guides the agent to the  
 176 expansion node. After following this abstract policy and reaching the expansion node, the agent  
 177 executes a low-level exploration policy (Section 3.2.6), and finally extracts a new node from the  
 178 resulting trajectory to add it to the skill graph (Section 3.2.7).

### 179 3.2.1 Identifying the Expansion Node

To select an expansion node, we need to quantify how much each node would contribute to graph  
 expansion: nodes that cause the graph to extend further should have a higher probability of being  
 selected for expansion. To capture this idea, we first construct an intrinsic reward function that is  
 higher the further away a state is from the graph. But rather than relying on distance functions, which  
 are usually unavailable in high-dimensional observation spaces, we use pseudocounts (Bellemare  
 et al., 2016): the further away a state is from the graph, the fewer times it has likely been visited.  
 We use Coin Flip Networks (CFN) (Lobel et al., 2023) to estimate pseudocounts and use it as our  
 intrinsic reward:

$$r_{\text{novelty}}(s) = \frac{1}{\sqrt{N_{\phi}(s)}},$$

180 where  $N_{\phi}(s)$  is the pseudocount of state  $s$  as approximated using CFN parameters  $\phi$ .

181 We now have an intrinsic reward function  $r_{\text{novelty}}$  that is higher along the frontiers of the graph,  
 182 where the agent has less experience. To turn this into an expansion probability, we ask: if the agent  
 183 executes a novelty-maximizing policy  $\pi_{\text{novelty}}$  starting from each graph node, how much intrinsic re-  
 184 ward can it expect to accumulate? Fortunately, this concept is well captured by the exploration value  
 185 function  $V_{\text{novelty}}$ . But, solely using novelty will result in graph expansion in all possible directions;  
 186 to forego exhaustive coverage and prioritize high reward regions, we add the intrinsic and extrinsic  
 187 rewards together while learning  $V_{\text{novelty}}$ . Concretely, we define the utility of expanding a node  $n$  as  
 188  $U(n) = \mathbb{E}_{s \sim n} [V_{\text{novelty}}(s)] = \mathbb{E}_{s \sim n} \left[ \sum_{t=1}^{\infty} R(s_t, a_t) + b \cdot r_{\text{novelty}}(s_t) \mid s_0 = s \right]$ ; we use  $b = 0.01$ .



189 Finally, we restrict the choice of expansion node to the set of nodes that are reachable from the  
 190 current state, i.e, the descendants of the nodes that map to the current state of the agent  $s_t$ . Therefore,  
 191 we define the probability of expanding a node  $n$  as follows:

$$n_{\text{expansion}} \sim \mathbb{P}_n = \frac{U(n)}{\sum_{n' \in \mathcal{D}(s_t)} U(n')} = \frac{\mathbb{E}_{s \sim n} [V_{\text{novelty}}(s)]}{\sum_{n' \in \mathcal{D}(s_t)} \mathbb{E}_{s' \sim n'} [V_{\text{novelty}}(s')]}, \forall n \in \mathcal{D}(s_t). \quad (4)$$

### 192 3.2.2 Estimating Edge Probabilities

193 Every edge in the skill graph has a weight  $p_{n_i \rightarrow n_j}$ , which represents the probability with which an  
 194 option  $o_{n_j}$  initiating from states in node  $n_i$  will successfully terminate in node  $n_j$ . As pointed out  
 195 earlier, this is equivalent to the expected initiation probability of option  $o_{n_j}$ , where the expectation  
 196 is taken with respect to states inside the node  $n_i$ . Following Bagaria et al. (2023), we interpret the  
 197 agent’s goal-conditioned value function (UVFA) (Schaul et al., 2015) as an initiation probability:

$$p_{n_i \rightarrow n_j} = \mathbb{E}_{s \sim n_i} [\mathcal{I}_{o_j}(s)] \approx \mathbb{E}_{s \sim n_i} [V_{n_j}(s)] := v_{n_i \rightarrow n_j}. \quad (5)$$

198 It is possible to infer the edge probabilities from the UVFA because the goal-conditioned reward  
 199 function that is used to train it is  $\mathcal{R}_o(s) = \beta_o(s)$ , i.e., it is +1 for states inside the option’s subgoal  
 200 region and 0 otherwise. However, naively using the UVFA prediction as edge probability can hinder  
 201 learning because:

- 202 1.  $v_{n_i \rightarrow n_j}$  for an edge can be low if traversing it is actually unlikely in a single option execution  
 203 (even under an optimal policy) *or* because the agent does not yet have enough data, and,
- 204 2. when we add new nodes to the graph, the probability of reaching them will initially be low. To  
 205 increase these edge probabilities, the agent must continually practice the options that target them,  
 206 but this is unlikely to happen owing to their low starting probabilities.

207 Therefore, to address these issues, we are *optimistic* about the value of the edge probabilities:

$$p_{n_i \rightarrow n_j}^+ := v_{n_i \rightarrow n_j} + \mathcal{U}(n_i, n_j) \approx v_{n_i \rightarrow n_j} + \frac{c}{\sqrt{N[n_i, n_j] + 1}}, \quad (6)$$

208 where  $\mathcal{U}(n_i, n_j)$  represents the uncertainty about our value prediction  $v_{n_i \rightarrow n_j}$ , which we obtain  
 209 using  $N[n_i, n_j]$  (Brafman & Tenenbholz, 2002; Strehl & Littman, 2008), the number of times the  
 210 agent has executed option  $o_{n_j}$  starting from states in node  $n_i$  ( $c \in [0, 1]$  is a hyperparameter).<sup>1</sup> When  
 211  $N[n_i, n_j]$  is low, the agent is optimistic about traversing that edge, and so it is encouraged to practice  
 212 the corresponding option policy. But, as the agent attempts to traverse edge  $e_{n_i \rightarrow n_j}$ , the count  
 213  $N[n_i, n_j]$  increases, and the optimistic bias decays; eventually, the edge probability approaches the  
 214 goal-conditioned value, which we expect to become more accurate with more data. By incorporating  
 215 an optimistic bias via the edge probabilities, we ensure that the policies computed by the planner  
 216 achieve a form of **exploration within the graph**, which prioritizes paths that reach the expansion  
 217 node while improving newly added skills.

### 218 3.2.3 Deciding when there is an edge between two nodes

219 Using Equation 6, we can determine the probability of traversing an edge in a single option exe-  
 220 cution. Naïvely, this would result in a dense, fully-connected graph, which would be problematic:  
 221 suppose there is an existing node  $n_1$  and we just added a new node  $n_2$  that is far away (even un-  
 222 der an optimal policy) from  $n_1$ . The agent’s initial desire to connect them would be high (because  
 223  $N[n_1, n_2]$  will be low). As a result, the agent will spend a lot of time realizing that  $n_1$  and  $n_2$  should  
 224 actually have a low probability of being directly connected to one another. Not only is this sample-  
 225 inefficient, it also results in graph construction *slowing down over time*: for every new node added

<sup>1</sup>Since the counts are over discrete states of the AMDP, we simply maintain tabular counts in a  $|\bar{S}| \times |\bar{S}|$  matrix.

to the graph, the agent must determine how the new node relates to every other existing node. To mitigate this issue, we use the following heuristics to determine when two nodes  $n_1$  and  $n_2$  should have an edge connecting them:

- **Observed multi-step transition.** When we observe a transition between two nodes  $n_1 \rightarrow n_2$  within  $H$  steps (the maximum horizon of an option), we connect them with an edge.
- **High goal-conditioned value.** We connect nodes  $n_1 \rightarrow n_2$  with an edge when the goal-conditioned value between them,  $\mathbb{E}_{s \sim n_1} [V_{n_2}(s)]$ , is high enough. To ensure that the resulting model is amenable to planning, we add edges only when we are *highly confident* that execution of one option will allow for the execution of another (Konidaris et al., 2018). We instantiate this principle via an adaptive threshold that is high to begin with, and tapers to a fixed hyperparameter  $\tau \in [0, 1]$  when more experience is gained:

$$e_{n_1 \rightarrow n_2} = \mathbb{I} \left[ \mathbb{E}_{s \sim n_1} [V_{n_2}(s)] - \frac{c}{\sqrt{\min \{N[n_1], N[n_2]\}}} > \tau \right], \quad (7)$$

where  $\mathbb{I}$  is the indicator function and  $N[n]$  is the number of times the agent has visited node  $n$ . The inverse square root term acts as an uncertainty penalty that diminishes as more data is collected.

### 3.2.4 Constructing the Abstract MDP

After the agent identifies the expansion node  $n_{\text{expansion}}$ , it needs to compute a plan that will guide it from its current state  $s_t$  to  $n_{\text{expansion}}$ . Our graph has two properties that will inform an appropriate choice of planning algorithm:

1. **Probabilistic transitions.** The skill graph has probabilistic edges because (a) the environment’s transition function and the agent’s option policies are stochastic, and (b) the option policies are being learned online, causing the abstract transition function to also appear stochastic and non-stationary (Nachum et al., 2018b; Levy et al., 2019; Bagaria et al., 2021a). We represent these dynamics with a transition matrix,  $\bar{T}$ , where each entry  $\bar{T}_{ij}$  is the probability that executing the option targeting node  $n_j$  from node  $n_i$  successfully leads to node  $n_j$ . When option execution is unsuccessful, we assume the agent “falls off” the graph, modeling it by assigning the remaining probability  $(1 - \bar{T}_{ij})$  to an absorbing failure state  $\varphi$ .
2. **Incorporating extrinsic reward.** Given two paths from  $s_t$  that both reach  $n_{\text{expansion}}$ , we would prefer one that collects extrinsic rewards along the way. This allows us to form a reward-respecting abstract model (Sutton et al., 2024), and prioritizes skill refinement in parts of the graph that are likely to result in higher reward. To incorporate this, we define a reward matrix,  $\bar{R}$ , where each entry  $\bar{R}_{ij}$  is the sum of extrinsic rewards collected during the execution of the option and an intrinsic reward that is high at the expansion node. Formally, we write:  $\bar{R}_{ij} = \kappa \bar{R}_{ij}^{\text{ext}} + \bar{R}_{ij}^{\text{int}}$ , where  $\kappa \in \mathbb{R}$  scales the extrinsic reward, and  $\bar{R}_{ij}^{\text{int}}$  equals 1 if node  $n_j$  is the expansion node, and 0 otherwise.

We also define an abstract discount factor,  $\bar{\gamma}$ , which reflects the probability of staying within the graph during an option execution. Putting these elements together, our AMDP is defined as:  $\bar{M} = (\bar{S}, \bar{A}, \bar{R}, \bar{T}, \bar{\gamma})$ , where  $\bar{S}$  and  $\bar{A}$  represent the set of nodes in the skill graph. Planning on this abstract MDP is a Stochastic Shortest Path (SSP) problem (Bertsekas & Tsitsiklis, 1991). We solve this SSP using value iteration (Bellman, 1966), which yields a deterministic abstract policy  $\bar{\pi}$ . This abstract policy maps each node (abstract state) to an option, guiding the agent toward the expansion node while balancing the collection of extrinsic rewards and exploration of uncertain edges. An illustrative example of this AMDP construction is provided in Figure 7 in Appendix C.

### 3.2.5 Navigating to the Expansion Node

The abstract policy  $\bar{\pi}$  tells the agent which options to execute to reach the expansion node. While following  $\bar{\pi}$ , if the agent finds itself in a state that is not in the graph ( $s_t \in \varphi$ ), it targets the



closest node in the graph:  $o_t = \operatorname{argmax}_o V_{\beta_o}(s_t), \forall o \in \mathcal{O}$ . Each option execution corresponds to rolling out the corresponding option policy  $\pi_o(s)$ , which is trained using recurrent deep Q-learning (Kapturowski et al., 2019) to maximize the pseudo reward function  $R_o(s) = \beta_o(s)$ . In practice, all option policies are parameterized using the same goal-conditioned policy  $\pi(s|g; \theta)$  (Schaul et al., 2015) and each option conditions the policy with goal states sampled from its own termination region (Bagaria et al., 2021a). To improve the sample-efficiency of learning effective option policies, we use hindsight experience replay (Andrychowicz et al., 2017).

### 3.2.6 Generating Behavior Outside the Graph

After reaching the expansion node, the IM-DSG agent attempts to extend the graph. It does so by executing the low-level exploration policy  $\pi_{\text{CFN}}$ , which is trained to maximize the weighted sum of intrinsic and extrinsic rewards:

$$R_{\text{CFN}}(s, a) = R(s, a) + b \cdot r_{\text{novelty}}(s),$$

where  $r_{\text{novelty}}(s) = 1/\sqrt{N_\phi(s)}$ ,  $b \in [0, 1]$  controls the exploration-exploitation trade-off (we use  $b = 0.01$ ) and  $N_\phi(s)$  is the CFN visitation count prediction for state  $s$ . We use recurrent deep Q-learning (Kapturowski et al., 2019) to learn  $V_{\text{novelty}}$  and  $\pi_{\text{novelty}}$ . The result of this policy execution is a trajectory  $\tau_{\text{novelty}}$ , which is likely to move outward, toward regions where the agent’s abstract model is more uncertain; this is a way of doing exploration *outside* the graph.

### 3.2.7 Adding a New Node to the Graph

Having generated a low-level exploration trajectory  $\tau_{\text{novelty}}$  in the frontier of the graph, the agent must now identify a new node to add to the graph. Again, we take inspiration from intrinsic motivation methods (Simsek & Barto, 2004): the agent identifies the state in  $\tau_{\text{novelty}}$  with the highest predicted novelty; if the novelty of that state is higher than the running mean of the novelty-based rewards, then we flag it as a new target and add it to the skill graph:

$$s_{\text{target}} = \operatorname{argmax}_{s \in \tau_{\text{novelty}}} [r_{\text{novelty}}(s)] \text{ if } \max_s [r_{\text{novelty}}(s)] > \mu_t + k\sigma_t, \text{ else } \varphi, \quad (8)$$

where  $\mu_t$  is the running mean and  $\sigma_t$  is the running standard deviation of  $r_{\text{novelty}}$ , computed incrementally over the lifetime of the agent (Welford, 1962);  $k \in \mathbb{R}^+$  is a hyperparameter (we use  $k = 1$ ). Equation 8 identifies a target state, which we need to convert to a termination classifier so that it can serve as a node in the skill graph. We could use domain knowledge (for example, the “cell” representations of Ecoffet et al. (2021)), pixel-wise equality (Veeriah et al., 2018), or a temporal distance based classifier (Savinov et al., 2018); we opt for simplicity and use an off-the-shelf template matching algorithm, which is a simple call to OpenCV’s `matchTemplate(s, g)` (Bradski, 2000):

$$n_{\text{new}} = \beta_{o_{\text{new}}}(s) = \text{matchTemplate}(s, s_{\text{target}}). \quad (9)$$

### 3.2.8 Summary of Graph Discovery Algorithm

Algorithm 1 summarizes IM-DSG. The agent maintains two value functions  $V_g, V_{\text{novelty}}$  and two policies  $\pi_g, \pi_{\text{novelty}}$ ; R2D2 (Kapturowski et al., 2019) is used to learn both sets of policies and value functions. The agent first chooses an expansion node using Equation 4. Which option is executed is then determined by solving the resulting Abstract MDP using value iteration. New nodes are added to the graph by executing  $\pi_{\text{novelty}}$  from the expansion node. Finally, graph edges are updated based on  $V_g$  and edge traversal counts  $N[i, j]$ . The choice of expansion node and subsequent low-level exploration encourage the skill graph to expand along directions that jointly maximize novelty and extrinsic reward.

**Algorithm 1** Intrinsic Motivation Directed Skill Graph (IM-DSG)

---

```

1: Initialize:
2: Initialize CFN to learn pseudocounts  $N_\phi$ , value function  $V_{\text{novelty}}$ , and policy  $\pi_{\text{novelty}}$ .
3: Initialize goal-conditioned R2D2 to learn value function  $V_g$  and policy  $\pi_g$ .
4: Initialize graph  $\mathcal{G}$  with null node  $\{\varphi\}$  and no edges.
5: while True do
6:   Sample expansion node using Equation 4.
7:   Construct AMDP using procedure described in Section 3.2.4.
8:   Solve AMDP using value iteration to get abstract policy  $\bar{\pi}$ .
9:   Follow abstract policy until agent reaches expansion node or episode terminates.
10:  if agent reached expansion node then
11:    Roll out exploration policy  $\pi_{\text{novelty}}$  to get trajectory  $\tau$ .
12:  end if
13:  Potentially add a new node from  $\tau$  to the graph using Equation 8.
14:  Update  $V_g$  and  $\pi_g$  using Hindsight Experience Replay (HER).
15:  Update  $\pi_{\text{novelty}}$  to maximize the sum of intrinsic and extrinsic rewards:  $R(s, a) + b \cdot r_{\text{novelty}}(s)$ .
16:  Update the edge traversal  $N[i, j]$  and node visitation counts  $N[i]$  using the full trajectory.
17:  Update the edge probabilities in the skill graph  $\mathcal{G}$  using  $V_g$  (Sections 3.2.2 and 3.2.3).
18: end while

```

---

## 4 Experiments

We evaluate the IM-DSG agent on five challenging exploration problems, all of which have image-based observations and discrete actions (details in Appendix A). We chose the MINIGRID (Chevalier-Boisvert et al., 2023) environments because they were identified as being challenging for existing RL algorithms by the survey of Colas et al. (2022). We chose SOKOBAN (Schrader, 2018) because the survey of Hamrick et al. (2020) identified it as being a problem in which both learning and planning are important for good performance. Finally, we chose VISUALTAXI and VISUALPINBALL because non-image versions of these domains have been used extensively to evaluate state abstraction and option discovery algorithms in isolation (for example, Lo et al., 2024; Konidaris & Barto, 2009; Dietterich, 2000).

### 4.1 Qualitative Evaluation

To build intuition, we first qualitatively evaluate the IM-DSG agent. Specifically, we want to understand what the skill-graph looks like, what parts of it the agent uses most frequently, and whether the skill graph aids more effective exploration.

#### 4.1.1 Visualizing the Skill Graph in MINIGRID

Figure 2 is a visualization of the skill graph at the end of training in MINIGRID-KEYCORRIDOR. To visualize each node in the graph, we randomly sample a state in which that termination condition was triggered; only the location and inventory is visualized here for simplicity. Each node is colored based on its probability of serving as the expansion node. Furthermore, the transparency  $\alpha_n$  of a node is its average abstract value:

$$\alpha_n = \max \left\{ \mathbb{E}_{n_e} \left[ \bar{V}_{n_e}(n) \right], 0.05 \right\},$$

where the expectation is over the node expansion probability distribution  $\mathbb{P}_n$  (Equation 4),  $\bar{V}_g(n)$  is the abstract value of node  $n$  when targeting node  $g$  (obtained via Value Iteration on the AMDP; Section 3.2.4). The transparency of an edge is the average transparency of the nodes that it connects.

**Understanding Figure 2.** Nodes in the top-right room have a higher expansion probability. This is because (a) navigating to the top-right room is novel as it requires unlocking the red door with

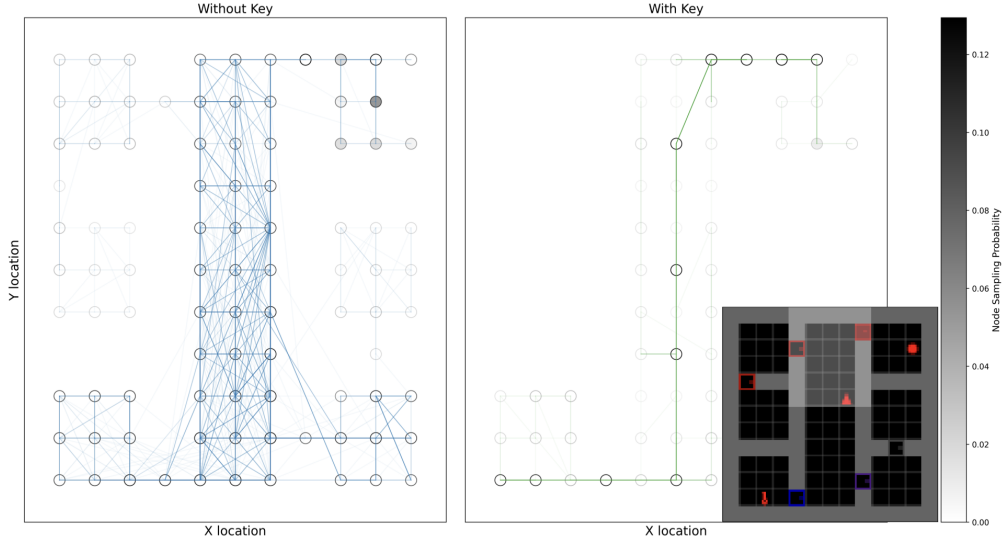


Figure 2: **Discovered skill graph in MINIGRID-KEYCORRIDOR.** The domain is shown in the bottom-right inset: it involves picking up a key in the bottom-left room, opening a locked door and then picking up a ball in the top-right room. The left subplot shows the subgraph where the agent does *not* have a key in its possession (other state variables are not shown for simplicity). The nodes and edges of the graph are shaded based on their expected abstract value  $\mathbb{E}[\bar{V}(n)]$ . Although we visualize the location and inventory, the agent itself only observes images. The graphs show that, when the agent lacks the key, it aims to explore down the corridor and in the rooms at the bottom of the maze, but when it does have the key, interesting options are on the path to the goal.

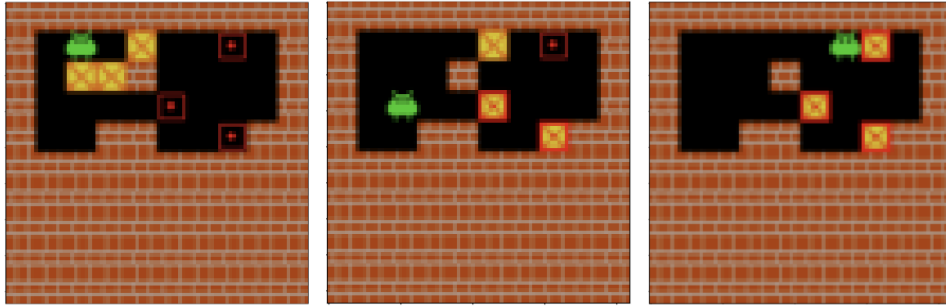


Figure 3: Example of hierarchical policy execution in SOKOBAN: (left) start state in which *no* box (boxes shown in yellow) is in its desired location (target locations are shown in red). (Middle) selected expansion node which has 2 out of 3 boxes placed in their desired locations; the agent uses planning to get to this state. (Right) After reaching the expansion node, the agent does local exploration using  $\pi_{\text{novelty}}$ , which reaches the task goal.

325 the red key, and (b) once the agent enters the top-right room, it can solve the remainder of the  
 326 task easily. These two facts result in high  $V_{\text{novelty}}$  and hence high expansion probabilities, which is  
 327 why the graph in the bottom-left and top-right rooms have high abstract value. Nodes in the central  
 328 corridor also have high abstract value because the agent must navigate through them for the majority  
 329 of its target expansion nodes. This demonstrates that the agent focuses on parts of the graph that are  
 330 either important for expansion by themselves, or those subgraphs that serve as *hubs*, i.e., they lead  
 331 to other sub-graphs that are important for expansion (Şimşek et al., 2005; Shah & Srivastava, 2024).

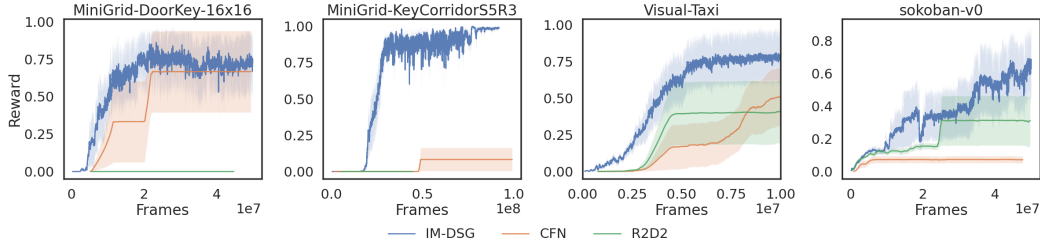


Figure 4: Comparing the performance of our proposed algorithm, IM-DSG (blue), novelty-based CFN (orange) and vanilla model-free R2D2 (green). Each curve represents the mean reward and the shaded area represents the standard deviation over 5 random seeds.

#### 4.1.2 Exploration in SOKOBAN

Figure 3 shows IM-DSG’s high-level policy when it encounters the task-goal for the first time: the expansion node is chosen to be a state that corresponds to placing two out of the three boxes in their correct locations. The agent plans with its learned options to reach this state. Subsequent novelty-driven exploration is more targeted as it only seeks to explore those states in which most of the task has already been solved, increasing the probability that random exploration would find the proverbial “needle in a haystack”.

### 4.2 Quantitative Evaluation

First, we seek to understand whether the IM-DSG agent can solve challenging exploration problems in a more sample-efficient manner than our baselines. Second, even though the graph construction is prioritized along directions where the agent expects to gather extrinsic reward, we want to understand whether the graph can still be used in a multi-task context to solve goal-reaching problems that were not encountered during training.

#### 4.2.1 Sample Efficient Exploration in Single Task MDPs

We compare the IM-DSG agent to the following baselines:

- **Flat model-free RL.** R2D2 is a distributed variant of DQN that uses a recurrent neural network to represent the value function (Kapturowski et al., 2019). Since the IM-DSG agent uses R2D2 to learn policies, we include it as our model-free baseline.
- **Novelty driven exploration.** We use CFN to explore outside the graph, and it is a state-of-the-art technique for novelty-driven RL; so, we include it as our exploration baseline.

Figure 4 shows our quantitative results: we find that the IM-DSG agent consistently outperforms R2D2 and CFN in our domains. This suggests that although novelty-maximizing exploration is in principle sufficient to solve these problems, the discovery and use of the model for planning leads to a more sample-efficient agent in these domains.

**Model-based planning ablation.** Figure 5 (left) shows that when we disable planning, the IM-DSG agent is unable to solve the MINIGRID-KEYCORRIDOR problem. This indicates that high-level planning with the learned skill graph is essential—simply combining CFN with HER is not enough to solve this task in the allotted training budget. This ablation also serves a rough comparison to algorithms like Skew-Fit (Pong et al., 2019) and Omega (Pitis et al., 2020) that combine HER with novelty-based exploration, but do not do model-based planning.

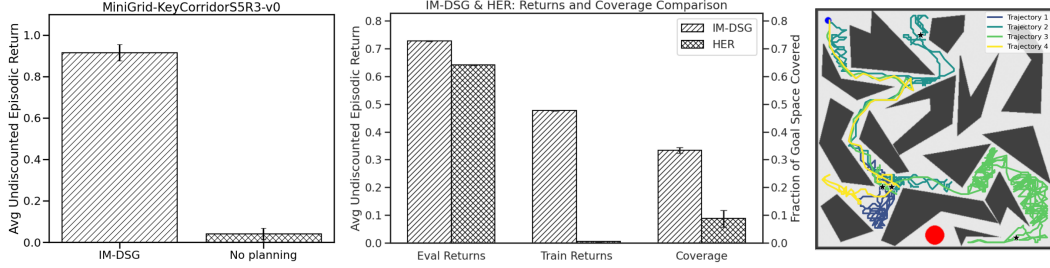


Figure 5: **Ablation and graph usefulness in a multi-task context.** (left) Ablating the use of planning in IM-DSG in MINIGRID-KEYCORRIDOR: bars represent the average per-episode reward during training. (middle) Evaluation performance in VISUALPINBALL on 10k randomly sampled goal states not used during training: comparison against HER in terms of test-time rewards, train-time rewards, and coverage, which is the fraction of the maze explored during training. All error bars represent standard deviation over 5 random seeds. (right) Solution trajectories found by the IM-DSG agent for 4 randomly sampled test-time goals (shown using black stars); the blue dot in the top-left of the maze is the start state, the big red ball in the bottom is the train-time goal. Note that, unlike in the original, this version of pinball uses images as observations.

#### 362 4.2.2 Generalization to New Goals in VISUALPINBALL

363 In this section, we test whether the skill graph discovered during exploration can be used to solve  
 364 other, unseen tasks at test-time. In particular, we train our agent to solve the VISUALPINBALL  
 365 problem, which requires guiding the ball to the big red circle at the bottom of the maze. After  
 366 training, we present the agent with 10k randomly sampled held-out goals and test its ability to reach  
 367 these goals. We compare our agent against Hindsight Experience Replay (HER; Andrychowicz  
 368 et al., 2017), which in principle can also generalize to unseen goals at test-time. We follow the  
 369 experimental setup of Sharma et al. (2020a); more details can be found in Appendix B.

370 Figure 5 shows that IM-DSG outperforms HER at test-time. To understand why, we compare their  
 371 train-time performances and find that HER is unable to explore the maze as extensively as IM-DSG,  
 372 as a result of which, it is unable to reach goals that are far away from the start-state at test-time.

### 373 5 Discussion and Conclusion

374 We introduced an algorithm that builds a graph abstraction of large MDPs with image-based ob-  
 375 servations. Nodes of this graph correspond to option termination regions, edges represent option  
 376 policies, and edge probabilities are determined using option initiation functions. The skill graph ex-  
 377 pands towards the frontier of the agent’s experience without assuming access to state-sampling or a  
 378 distance metric; the acquired graph enables planning in increasingly large portions of the state-space.

379 Our method suffers from at least the following limitations:

- 380 • Each option drives the value of all state variables to a certain range of values. In more complex  
 381 environments, it may be unnecessary, or even impossible, to control all state variables at once.
- 382 • After reaching the expansion node, the IM-DSG agent executes the low-level exploration policy  
 383 (CFN) until the end of the episode; more fine-grained control over when to enter and exit low-level  
 384 exploration (Pislar et al., 2022) would generalize our algorithm to continuing problems without  
 385 episode resets.

386 Despite these shortcomings, IM-DSG takes a small step towards the eventual goal of building RL  
 387 agents that continually discover their own state and action abstractions (Gershman, 2017; Konidaris,  
 388 2019), and use them for effective exploration (Gopnik, 2020) and planning (Sutton et al., 2024).

## References

- Cameron Allen, Neev Parikh, Omer Gottesman, and George Konidaris. Learning markov state abstractions for deep reinforcement learning. *Advances in Neural Information Processing Systems*, 34:8229–8241, 2021.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In *Advances in Neural Information Processing Systems*, pp. 5048–5058, 2017.
- Pierre-Luc Bacon. *Temporal Representation Learning*. PhD thesis, McGill University Libraries, 2018.
- Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- Akhil Bagaria and Tom Schaul. Scaling goal-based exploration via pruning proto-goals. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJ-CAI 2023, 19th-25th August 2023, Macao, SAR, China*, pp. 3451–3460. ijcai.org, 2023. DOI: 10.24963/IJCAI.2023/384. URL <https://doi.org/10.24963/ijcai.2023/384>.
- Akhil Bagaria, Jason Senthil, Matthew Slivinski, and George Konidaris. Robustly learning composable options in deep reinforcement learning. In *30th International Joint Conference on Artificial Intelligence*, 2021a.
- Akhil Bagaria, Jason K. Senthil, and George Konidaris. Skill discovery for exploration and planning using deep skill graphs. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 521–531. PMLR, 2021b. URL <http://proceedings.mlr.press/v139/bagaria21a.html>.
- Akhil Bagaria, Ben M Abbatematteo, Omer Gottesman, Matt Corsaro, Sreehari Rammohan, and George Konidaris. Effectively learning initiation sets in hierarchical reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=4JCVw8oMlf>.
- Andrew G. Barto. Intrinsic motivation and reinforcement learning. In Gianluca Baldassarre and Marco Mirolli (eds.), *Intrinsically Motivated Learning in Natural and Artificial Systems*, pp. 17–47. Springer, 2013. DOI: 10.1007/978-3-642-32375-1\_2. URL [https://doi.org/10.1007/978-3-642-32375-1\\_2](https://doi.org/10.1007/978-3-642-32375-1_2).
- Andrew G Barto and Sridhar Mahadevan. Recent advances in hierarchical reinforcement learning. *Discrete event dynamic systems*, 13(1-2):41–77, 2003.
- Kate Baumli, David Warde-Farley, Steven Hansen, and Volodymyr Mnih. Relative variational intrinsic control. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021*, pp. 6732–6740, 2021.
- Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems*, pp. 1471–1479, 2016.
- Richard Bellman. Dynamic programming. *Science*, 153(3731):34–37, 1966.
- Dimitri P Bertsekas and John N Tsitsiklis. An analysis of stochastic shortest path problems. *Mathematics of Operations Research*, 16(3):580–595, 1991.
- G. Bradski. The OpenCV Library. *Dr. Dobb’s Journal of Software Tools*, 2000.



- 432 Ronen I Brafman and Moshe Tennenholtz. R-Max - a general polynomial time algorithm for near-  
433 optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct):213–231, 2002.
- 434 Víctor Campos Camúñez, Alex Trott, Caiming Xiong, Richard Socher, Xavier Giró Nieto, and Jordi  
435 Torres Viñals. Explore, discover and learn: unsupervised discovery of state-covering skills. In  
436 *Proceedings of the 37th International Conference on Machine Learning, ICML*, volume 119, pp.  
437 1317–1327, 2020.
- 438 Nuttapon Chentanez, Andrew G Barto, and Satinder P Singh. Intrinsically motivated reinforcement  
439 learning. In *Advances in neural information processing systems*, pp. 1281–1288, 2005.
- 440 Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems,  
441 Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld:  
442 Modular & customizable reinforcement learning environments for goal-oriented tasks. *CoRR*,  
443 abs/2306.13831, 2023.
- 444 Jongwook Choi, Archit Sharma, Honglak Lee, Sergey Levine, and Shixiang Shane Gu. Variational  
445 empowerment as representation learning for goal-conditioned reinforcement learning. In *Interna-*  
446 *tional Conference on Machine Learning*, pp. 1953–1963. PMLR, 2021.
- 447 Fan RK Chung and Fan Chung Graham. *Spectral graph theory*. Number 92. American Mathematical  
448 Soc., 1997.
- 449 Cédric Colas, Tristan Karch, Olivier Sigaud, and Pierre-Yves Oudeyer. Autotelic agents with intrin-  
450 sically motivated goal-conditioned reinforcement learning: a short survey. *Journal of Artificial*  
451 *Intelligence Research*, 74:1159–1199, 2022.
- 452 Peter Dayan and Geoffrey E Hinton. Feudal reinforcement learning. In *Advances in neural infor-*  
453 *mation processing systems*, pp. 271–278, 1993.
- 454 Thomas G Dietterich. Hierarchical reinforcement learning with the maxq value function decompo-  
455 sition. *Journal of Artificial Intelligence Research*, 13:227–303, 2000.
- 456 Carlos Diuk, Andre Cohen, and Michael L. Littman. An object-oriented representation for ef-  
457 ficient reinforcement learning. In *Proceedings of the 25th International Conference on Ma-*  
458 *chine Learning, ICML ’08*, pp. 240–247, New York, NY, USA, 2008. Association for Com-  
459 puting Machinery. ISBN 9781605582054. DOI: 10.1145/1390156.1390187. URL <https://doi.org/10.1145/1390156.1390187>.  
460 <https://doi.org/10.1145/1390156.1390187>.
- 461 Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O Stanley, and Jeff Clune. First return, then  
462 explore. *Nature*, 590(7847):580–586, 2021.
- 463 Ben Eysenbach, Russ R Salakhutdinov, and Sergey Levine. Search on the replay buffer: Bridging  
464 planning and reinforcement learning. In *Advances in Neural Information Processing Systems 32*,  
465 pp. 15246–15257. 2019a.
- 466 Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need:  
467 Learning skills without a reward function. In *International Conference on Learning Representa-*  
468 *tions*, 2019b. URL <https://openreview.net/forum?id=SJx63jRqFm>.
- 469 Samuel J Gershman. On the blessing of abstraction. *SAGE Publications Sage UK: London, England*,  
470 2017.
- 471 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- 472 Alison Gopnik. Childhood as a solution to explore–exploit tensions. *Philosophical Transactions of*  
473 *the Royal Society B*, 375(1803):20190502, 2020.
- 474 Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *arXiv*  
475 *preprint arXiv:1611.07507*, 2016.

- 476 Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James  
477 Davidson. Learning latent dynamics for planning from pixels. *arXiv preprint arXiv:1811.04551*,  
478 2018.
- 479 Danijar Hafner, Kuang-Huei Lee, Ian Fischer, and Pieter Abbeel. Deep hierarchical planning from  
480 pixels. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.),  
481 *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Informa-*  
482 *tion Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December*  
483 *9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/a766f56d2da42cae20b5652970ec04ef-Abstract-Conference.html)  
484 [a766f56d2da42cae20b5652970ec04ef-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/a766f56d2da42cae20b5652970ec04ef-Abstract-Conference.html).
- 485 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy P. Lillicrap. Mastering diverse domains  
486 through world models. *CoRR*, abs/2301.04104, 2023. DOI: 10.48550/ARXIV.2301.04104. URL  
487 <https://doi.org/10.48550/arXiv.2301.04104>.
- 488 Jessica B Hamrick, Abram L Friesen, Feryal Behbahani, Arthur Guez, Fabio Viola, Sims Wither-  
489 spoon, Thomas Anthony, Lars Buesing, Petar Veličković, and Théophane Weber. On the role of  
490 planning in model-based deep reinforcement learning. *arXiv preprint arXiv:2011.04021*, 2020.
- 491 Steven Hansen, Guillaume Desjardins, Kate Baumli, David Warde-Farley, Nicolas Heess, Simon  
492 Osindero, and Volodymyr Mnih. Entropic desired dynamics for intrinsic control. *Advances in*  
493 *Neural Information Processing Systems*, 34:11436–11448, 2021.
- 494 Zhiao Huang, Fangchen Liu, and Hao Su. Mapping state space using landmarks for uni-  
495 versal goal reaching. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Flo-  
496 rence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural In-*  
497 *formation Processing Systems 32: Annual Conference on Neural Information Process-*  
498 *ing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp.  
499 1940–1950, 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/](https://proceedings.neurips.cc/paper/2019/hash/3b712de48137572f3849aabd5666a4e3-Abstract.html)  
500 [3b712de48137572f3849aabd5666a4e3-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/3b712de48137572f3849aabd5666a4e3-Abstract.html).
- 501 Yuu Jinnai, Jee Won Park, David Abel, and George Konidaris. Discovering options for exploration  
502 by minimizing cover time. In *Proceedings of the 36th International Conference on Machine*  
503 *Learning*. PMLR, 2019. URL [http://proceedings.mlr.press/v97/jinnai19b.](http://proceedings.mlr.press/v97/jinnai19b.html)  
504 [html](http://proceedings.mlr.press/v97/jinnai19b.html).
- 505 Yuu Jinnai, Jee Won Park, Marlos C. Machado, and George Konidaris. Exploration in reinforcement  
506 learning with deep covering options. In *International Conference on Learning Representations*,  
507 2020. URL <https://openreview.net/forum?id=SkeIyaVtwB>.
- 508 Leslie Pack Kaelbling. Learning to achieve goals. In *Proceedings of the 13th International Joint*  
509 *Conference on Artificial Intelligence (IJCAI)*, pp. 1094–1099. Citeseer, 1993.
- 510 Steven Kapturowski, Georg Ostrovski, Will Dabney, John Quan, and Remi Munos. Recurrent ex-  
511 perience replay in distributed reinforcement learning. In *International Conference on Learning*  
512 *Representations*, 2019. URL <https://openreview.net/forum?id=r1lyTjAqYX>.
- 513 Steven Kapturowski, Víctor Campos, Ray Jiang, Nemanja Rakićević, Hado van Hasselt, Charles  
514 Blundell, and Adrià Puigdomènech Badia. Human-level atari 200x faster. *arXiv preprint*  
515 *arXiv:2209.07550*, 2022.
- 516 Martin Klissarov and Marlos C Machado. Deep laplacian-based options for temporally-extended  
517 exploration. *arXiv preprint arXiv:2301.11181*, 2023.
- 518 G. Konidaris. *Autonomous Robot Skill Acquisition*. PhD thesis, University of Massachusetts at  
519 Amherst, 2011.
- 520 George Konidaris. On the necessity of abstraction. *Current Opinion in Behavioral Sciences*, 29:  
521 1–7, 2019.

- George Konidaris and Andrew Barto. Skill discovery in continuous reinforcement learning domains using skill chaining. In *Advances in Neural Information Processing Systems*, pp. 1015–1023, 2009.
- George Konidaris, Leslie Pack Kaelbling, and Tomas Lozano-Perez. From skills to symbols: Learning symbolic representations for abstract high-level planning. *Journal of Artificial Intelligence Research*, 61:215–289, 2018.
- Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- Steven M LaValle. Rapidly-exploring random trees: A new tool for path planning. 1998.
- Seungjae Lee, Jigang Kim, Inkyu Jang, and H. Jin Kim. DHRL: A graph-based approach for long-horizon and sparse hierarchical reinforcement learning. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/58b286aea34a91a3d33e58af0586fa40-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/58b286aea34a91a3d33e58af0586fa40-Abstract-Conference.html).
- Andrew Levy, George Konidaris, Robert Platt, and Kate Saenko. Hierarchical reinforcement learning with hindsight. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=ryzECoAcY7>.
- Stephen R Lindemann and Steven M LaValle. Incrementally reducing dispersion by increasing voronoi bias in RRTs. In *IEEE International Conference on Robotics and Automation*, volume 4, pp. 3251–3257, 2004.
- Chunlok Lo, Gabor Mihucz, Adam White, Farzane Aminmansour, and Martha White. Goal-space planning with subgoal models. *CoRR*, abs/2206.02902, 2022. DOI: 10.48550/ARXIV.2206.02902. URL <https://doi.org/10.48550/arXiv.2206.02902>.
- Chunlok Lo, Kevin Roice, Parham Mohammad Panahi, Scott M Jordan, Adam White, Gabor Mihucz, Farzane Aminmansour, and Martha White. Goal-space planning with subgoal models. *Journal of Machine Learning Research*, 25(330):1–57, 2024.
- Sam Lobel, Akhil Bagaria, and George Konidaris. Flipping coins to estimate pseudocounts for exploration in reinforcement learning. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 22594–22613. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/lobel23a.html>.
- M Chlodovskis Machado. *Efficient Exploration in Reinforcement Learning through Time-Based Representations*. PhD thesis, PhD thesis, University of Alberta, Canada, 2019.
- Marlos C. Machado, Marc G. Bellemare, and Michael Bowling. A laplacian framework for option discovery in reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, pp. 2295–2304, 2017.
- Marlos C Machado, Marc G Bellemare, Erik Talvitie, Joel Veness, Matthew Hausknecht, and Michael Bowling. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *Journal of Artificial Intelligence Research*, 61:523–562, 2018.
- Marlos C. Machado, Andre Barreto, and Doina Precup. Temporal abstraction in reinforcement learning with the successor representation. *Journal of Machine Learning Research (JMLR)*, 24(80):1–69, 2023.

- 566 Sridhar Mahadevan and Mauro Maggioni. Proto-value functions: A laplacian framework for learn-  
567 ing representation and control in markov decision processes. *Journal of Machine Learning Re-*  
568 *search*, 8(Oct):2169–2231, 2007.
- 569 Timothy A Mann, Shie Mannor, and Doina Precup. Approximate value iteration with temporally  
570 extended actions. *Journal of Artificial Intelligence Research*, 53:375–438, 2015.
- 571 Willie McClinton, Andrew Levy, and George Konidaris. HAC explore: Accelerating explo-  
572 ration with hierarchical reinforcement learning. *CoRR*, abs/2108.05872, 2021. URL <https://arxiv.org/abs/2108.05872>.  
573 <https://arxiv.org/abs/2108.05872>.
- 574 Elizabeth Amy McGovern. *Autonomous discovery of temporal abstractions from interaction with an*  
575 *environment*. PhD thesis, University of Massachusetts at Amherst, 2002.
- 576 Russell Mendonca, Oleh Rybkin, Kostas Daniilidis, Danijar Hafner, and Deepak Pathak. Dis-  
577 covering and achieving goals via world models. In Marc’Aurelio Ranzato, Alina Beygelz-  
578 imer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances*  
579 *in Neural Information Processing Systems 34: Annual Conference on Neural Informa-*  
580 *tion Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 24379–  
581 24391, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/](https://proceedings.neurips.cc/paper/2021/hash/cc4af25fa9d2d5c953496579b75f6f6c-Abstract.html)  
582 [cc4af25fa9d2d5c953496579b75f6f6c-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/cc4af25fa9d2d5c953496579b75f6f6c-Abstract.html).
- 583 Shakir Mohamed and Danilo Jimenez Rezende. Variational information maximisation for intrinsi-  
584 cally motivated reinforcement learning. *Advances in neural information processing systems*, 28,  
585 2015.
- 586 Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Near-optimal representation learning  
587 for hierarchical reinforcement learning. *arXiv preprint arXiv:1810.01257*, 2018a.
- 588 Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical  
589 reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 3303–3313,  
590 2018b.
- 591 S. Nasiriany, V. Pong, S. Lin, and S. Levine. Planning with goal-conditioned policies. *Advances in*  
592 *Neural Information Processing Systems (NeurIPS)*, 2019.
- 593 Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations:  
594 Theory and application to reward shaping. In *International Conference on Machine Learning*,  
595 volume 99, pp. 278–287, 1999.
- 596 Junhyuk Oh, Yijie Guo, Satinder Singh, and Honglak Lee. Self-imitation learning. In *International*  
597 *conference on machine learning*, pp. 3878–3887. PMLR, 2018.
- 598 Miruna Pislari, David Szepesvari, Georg Ostrovski, Diana L. Borsa, and Tom Schaul. When should  
599 agents explore? In *The Tenth International Conference on Learning Representations, ICLR 2022,*  
600 *Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL [https://openreview.net/](https://openreview.net/forum?id=dEwfxtl4bca)  
601 [forum?id=dEwfxtl4bca](https://openreview.net/forum?id=dEwfxtl4bca).
- 602 Silviu Pitis, Harris Chan, Stephen Zhao, Bradly C. Stadie, and Jimmy Ba. Maximum entropy gain  
603 exploration for long horizon multi-goal reinforcement learning. In *Proceedings of the 37th In-*  
604 *ternational Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, vol-  
605 *ume 119 of Proceedings of Machine Learning Research*, pp. 7750–7761. PMLR, 2020. URL  
606 <http://proceedings.mlr.press/v119/pitis20a.html>.
- 607 Vitchyr H Pong, Murtaza Dalal, Steven Lin, Ashvin Nair, Shikhar Bahl, and Sergey Levine. Skew-  
608 fit: State-covering self-supervised reinforcement learning. *Proceedings of the 37th International*  
609 *Conference on Machine Learning, ICML, 2019*.

- 610 Doina Precup. *Temporal abstraction in reinforcement learning*. PhD thesis, University of Mas-  
611 sachusetts at Amherst, 2001.
- 612 Balaraman Ravindran. *An algebraic approach to abstraction in reinforcement learning*. PhD thesis,  
613 University of Massachusetts at Amherst, 2004.
- 614 Rafael Rodriguez-Sanchez and George Konidaris. Learning abstract world models for value-  
615 preserving planning with options. *RLJ*, 4:1733–1758, 2024.
- 616 Nikolay Savinov, Alexey Dosovitskiy, and Vladlen Koltun. Semi-parametric topological memory  
617 for navigation. In *International Conference on Learning Representations*, 2018.
- 618 Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approxima-  
619 tors. In *International conference on machine learning*, pp. 1312–1320, 2015.
- 620 Jürgen Schmidhuber. Formal theory of creativity, fun, and intrinsic motivation (1990-2010). *IEEE*  
621 *Trans. Auton. Ment. Dev.*, 2(3):230–247, 2010. DOI: 10.1109/TAMD.2010.2056368. URL  
622 <https://doi.org/10.1109/TAMD.2010.2056368>.
- 623 Max-Philipp B. Schrader. gym-sokoban. [https://github.com/mpSchrader/](https://github.com/mpSchrader/gym-sokoban)  
624 [gym-sokoban](https://github.com/mpSchrader/gym-sokoban), 2018.
- 625 Julian Schrittwieser, Thomas Hubert, Amol Mandhane, Mohammadamin Barekatain, Ioannis  
626 Antonoglou, and David Silver. Online and offline reinforcement learning by planning with a  
627 learned model. *Advances in Neural Information Processing Systems*, 34:27580–27591, 2021.
- 628 Naman Shah and Siddharth Srivastava. Hierarchical planning and learning for robots in stochastic  
629 settings using zero-shot option invention. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam  
630 Natarajan (eds.), *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-*  
631 *Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth*  
632 *Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024,*  
633 *Vancouver, Canada*, pp. 10358–10367. AAAI Press, 2024. DOI: 10.1609/AAAI.V38I9.28903.  
634 URL <https://doi.org/10.1609/aaai.v38i9.28903>.
- 635 Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-  
636 aware unsupervised discovery of skills. In *International Conference on Learning Representations*  
637 (*ICLR*), 2020a.
- 638 Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware  
639 unsupervised discovery of skills. (arXiv:1907.01657), February 2020b. URL [http://arxiv.](http://arxiv.org/abs/1907.01657)  
640 [org/abs/1907.01657](http://arxiv.org/abs/1907.01657). arXiv:1907.01657 [cs, stat].
- 641 Özgür Simsek and Andrew G. Barto. Using relative novelty to identify useful temporal abstrac-  
642 tions in reinforcement learning. In Carla E. Brodley (ed.), *Machine Learning, Proceedings*  
643 *of the Twenty-first International Conference (ICML 2004), Banff, Alberta, Canada, July 4-8,*  
644 *2004*, volume 69 of *ACM International Conference Proceeding Series*. ACM, 2004. DOI:  
645 10.1145/1015330.1015353. URL <https://doi.org/10.1145/1015330.1015353>.
- 646 Özgür Şimşek, Alicia P Wolfe, and Andrew G Barto. Identifying useful subgoals in reinforcement  
647 learning by local graph partitioning. In *Proceedings of the 22nd international conference on*  
648 *Machine learning*, pp. 816–823. ACM, 2005.
- 649 Alexander L Strehl and Michael L Littman. An analysis of model-based interval estimation for  
650 markov decision processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.
- 651 Richard S. Sutton. Integrated architectures for learning, planning, and reacting based on  
652 approximating dynamic programming. In Bruce W. Porter and Raymond J. Mooney  
653 (eds.), *Machine Learning, Proceedings of the Seventh International Conference on Ma-*  
654 *chine Learning, Austin, Texas, USA, June 21-23, 1990*, pp. 216–224. Morgan Kaufmann,



1990. DOI: 10.1016/B978-1-55860-141-3.50030-4. URL <https://doi.org/10.1016/b978-1-55860-141-3.50030-4>.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Richard S Sutton, Joseph Modayil, Michael Delp, Thomas Degris, Patrick M Pilarski, Adam White, and Doina Precup. Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*, pp. 761–768. International Foundation for Autonomous Agents and Multiagent Systems, 2011.
- Richard S Sutton, Michael Bowling, and Patrick M Pilarski. The Alberta plan for AI research. *arXiv preprint arXiv:2208.11173*, 2022.
- Richard S. Sutton, Marlos C. Machado, G. Zacharias Holland, David Szepesvari, Finbarr Timbers, Brian Tanner, and Adam White. Reward-respecting subtasks for model-based reinforcement learning (abstract reprint). In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (eds.), *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2024, February 20-27, 2024, Vancouver, Canada*, pp. 22713. AAAI Press, 2024. DOI: 10.1609/AAAI.V38I20.30613. URL <https://doi.org/10.1609/aaai.v38i20.30613>.
- R.S. Sutton, , D. Precup, and S. Singh. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1):181–211, 1999.
- Adrien Ali Taïga, William Fedus, Marlos C Machado, Aaron Courville, and Marc G Bellemare. Benchmarking bonus-based exploration methods on the arcade learning environment. *arXiv preprint arXiv:1908.02388*, 2019.
- Erik Talvitie. Self-correcting models for model-based reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.
- Vivek Veeriah, Junhyuk Oh, and Satinder Singh. Many-goals reinforcement learning. *arXiv preprint arXiv:1806.09605*, 2018.
- Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3540–3549, 2017.
- B. P. Welford. Note on a method for calculating corrected sums of squares and products. *Technometrics*, 4:419–420, 1962. URL <https://api.semanticscholar.org/CorpusID:120126049>.
- Tom Zahavy, Brendan O’Donoghue, Andre Barreto, Sebastian Flennerhag, Volodymyr Mnih, and Satinder Singh. Discovering diverse nearly optimal policies with successor features. In *ICML Workshop on Unsupervised Reinforcement Learning*, 2021.
- Lunjun Zhang, Ge Yang, and Bradly C. Stadie. World model as a graph: Learning latent landmarks for planning. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 12611–12620. PMLR, 2021. URL <http://proceedings.mlr.press/v139/zhang21x.html>.



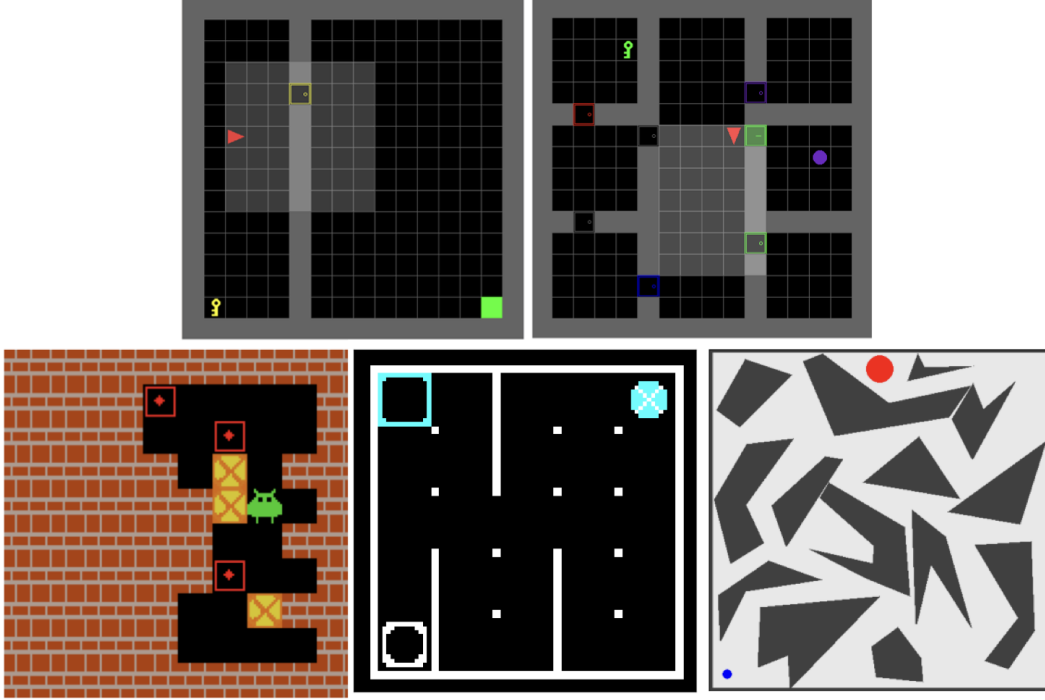
697 **A Domain Descriptions**

Figure 6: Domains used to test the IM-DSG agent. From top-left in left to right order: MINIGRID-DOORKEY, MINIGRID-KEYCORRIDOR, SOKOBAN, VISUALTAXI, and VISUALPINBALL.

- 698 1. MINIGRID-DOORKEY: To reach the goal, the agent must first pick up a key and then use that  
 699 key to unlock a door. There is no intermediate reward for picking up the key, only a sparse  
 700 terminating reward for reaching the goal. Each episode lasts a maximum of 200 steps.
- 701 2. MINIGRID-KEYCORRIDOR: This domain also has a key and a locked door, but additionally has  
 702 other doors that can be open and closed. The goal is to pick up the purple ball that is placed in  
 703 locked room (center-right room in Figure 6(c)). Each episode lasts a maximum of 1000 steps.
- 704 3. SOKOBAN: In this classic puzzle, the agent must place 3 boxes into their target locations (red  
 705 border squares in Figure 6). The environment provides intermediate rewards for putting each box  
 706 into its place and a terminating reward for correctly placing the final box. Each episode lasts a  
 707 maximum of 200 steps.
- 708 4. VISUALTAXI: this is an image-based version of the classic Taxi problem (Dietterich, 2000). A  
 709 passenger awaits in one of four depots and must be dropped off at a destination depot. Success-  
 710 ful completion of the full task yeilds a sparse terminating reward of +1. Each episode lasts a  
 711 maximum of 50 steps.
- 712 5. VISUALPINBALL: image-based version of the under-actuated Pinball domain (Konidaris &  
 713 Barto, 2009; Bacon et al., 2017) in which the agent provides acceleration to the ball in one  
 714 of 4 directions (or no-op), causing the velocity of the ball to change over time. We pick the  
 715 hardest configuration from Konidaris & Barto (2009) during training, and sample goal locations  
 716 at random during testing. Test-time goals are communicated to the agent as images of the pinball  
 717 in the desired location. Each episode lasts a maximum of 1000 steps.

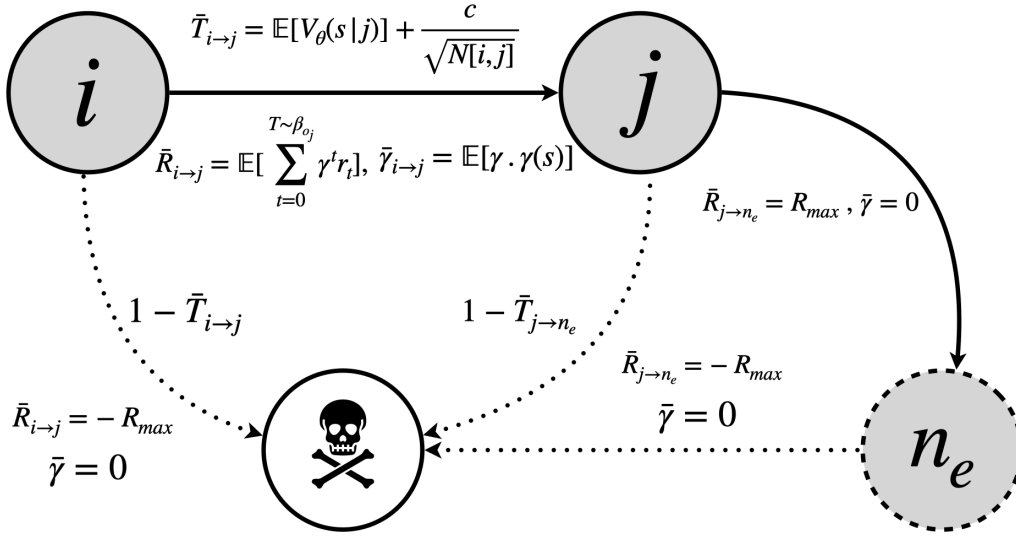


Figure 7: Summary illustration of a toy abstract MDP (AMDP) with 3 nodes.  $i$  and  $j$  are two ordinary nodes,  $n_e$  is the expansion node and the skull node shows a fictitious node that represents “falling off the graph”. The transition probabilities  $\bar{T}_{i \rightarrow j}$  are obtained using the UVFA  $V_\theta$  and edge-traversal counts  $N[i, j]$ ; the reward  $\bar{R}_{i \rightarrow j}$  and discount models  $\bar{\gamma}$  are obtained via monte carlo estimation from option executions.

## B VISUALPINBALL Experiment Details

Both the IM-DSG agent and the Hindsight Experience Replay (HER) agent get the same training budget of 50M frames. During evaluation, both algorithms are tested on the same set of 10k goals and are given a budget of 10M frames to achieve them. Each episode lasts for at most for 1k steps. Error bars are computed over 5 training runs of the respective agents.

To compute the coverage metric shown in Figure 5 (right), we first need to compute the maximum possible coverage in the maze. This is done by first discretizing the  $(x, y)$  plane into  $100 \times 100$  squares, and then subtracting the approximate area occupied by the obstacles in the maze. To compute the coverage achieved by the two agents, we consider what fraction of the free space is occupied by all the states visited by the two agents.

**Test-time modifications for VISUALPINBALL.** During training, IM-DSG picks the expansion node based on  $V_{\text{novelty}}$ . At test-time, this would not be a good strategy because there is no reason that  $V_{\text{novelty}}$  would guide the agent to the test-time goal. So instead, we pick the expansion node to be the one that is closest to the test-time goal  $g_{\text{test}}$ :

$$n_{\text{expansion}}^{\text{test}} = \underset{n}{\operatorname{argmax}} \mathbb{E}_{s \sim n} [V_{g_{\text{test}}}(s)] \quad (10)$$

## C AMDP Construction

Figure 7 illustrates a 3 node toy AMDP:  $n_e$  is the expansion node,  $i$  and  $j$  are intermediate nodes and the skull represents the null node  $\varphi$ . Nodes of the graph are abstract states and outgoing edges from a node are available options. The abstract model, which is a combination of  $\bar{T}$ ,  $\bar{R}$ ,  $\bar{\gamma}$ , is estimated using the agent’s goal-conditioned value function  $V_g$  and Monte Carlo estimation respectively.

Hyperparameter	Value	Tuned
Learning rate	<b>0.0001</b> , 0.00001	Yes
Trace length	40	No
Sequence period	20	No
Batch size	32	No
Burn-in length	0	No
Max replay size	500000	No
Target update period	<b>600</b> , 1200	Yes
Discount factor	0.997	No
Number of actors	32	No

Table 1: R2D2 Hyperparameters.

Hyperparameter	Value	Tuned
$\tau$ from Eq 5	0.7, <b>0.8</b> , 0.9	Yes
$c$ from Eq 6	0.1, <b>0.2</b>	Yes
$k$ in Eq 8	<b>1</b> , 2	Yes

Table 2: IM-DSG Specific Hyperparameters.

## D UVFA Training

We represent the agent’s goal-conditioned value function using the following neural network architecture: a convolutional neural network (CNN; Goodfellow et al., 2016) torso encodes the image representing the current state  $s$  into  $\Phi_1(s)$ ; another CNN encodes the goal  $g$  into  $\Phi_2(g)$ , where  $g$  is sampled from an option’s effect set.  $\Phi_1(s)$  is input into an LSTM (Goodfellow et al., 2016); the output of the LSTM,  $y_t$  is concatenated with  $\Phi_2(g)$  and then passed through a multi-layered perceptron (MLP), which outputs the Q-value for each action in the MDP’s action space. The goal encoder is identical to the state encoder, whose exact architecture, along with that of the LSTM and MLP, are taken without modification from Kapturowski et al. (2019). We follow Bagaria & Schaul (2023) and use novelty to pick which of the achieved goals in a trajectory should be replayed in hindsight.

## E Hyperparameters

Table 1 lists the hyperparameters for R2D2, Table 2 lists those hyperparameter settings that are specific to IM-DSG, and finally, Table 3 lists those IM-DSG hyperparameters that were environment specific.

Hyperparameter	DoorKey	KeyCorridor	VisualTaxi	VisualPinball	Sokoban
Option Horizon $H$	200	400	50	200	50
Descendant threshold	0	0	0.3	0	0

Table 3: IM-DSG Hyperparameter settings for different domains