
SkiLD: Unsupervised Skill Discovery Guided by Local Dependencies

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Unsupervised skill discovery carries the promise that an intelligent agent can learn
2 reusable skills through autonomous, reward-free environment interaction. Existing
3 unsupervised skill discovery methods learn skills by encouraging distinguishable
4 behaviors that cover diverse states. However, in complex environments with
5 many state factors (e.g., household environments with many objects), learning
6 skills that cover all possible states is impossible, and naively encouraging state
7 diversity often leads to simple skills that are not ideal for solving downstream
8 tasks. This work introduces Skill Discovery from Local Dependencies (SkiLD),
9 which leverages state factorization as a natural inductive bias to guide the skill
10 learning process. The key intuition guiding SkiLD is that skills that induce **diverse**
11 **interactions** between state factors are often more valuable for solving downstream
12 tasks. To this end, SkiLD develops a novel skill learning objective that explicitly
13 encourages the mastering of skills that effectively induce different interactions
14 within an environment. We evaluate SkiLD in several domains with challenging,
15 long-horizon sparse reward tasks including a realistic simulated household robot
16 domain, where SkiLD successfully learns skills with clear semantic meaning and
17 shows superior performance compared to existing unsupervised reinforcement
18 learning methods that only maximize state coverage.

19 1 Introduction

20 Reinforcement learning (RL) achieves impressive successes when solving decision-making problems
21 with well-defined reward functions [62, 19, 31]. However, designing this reward function is often
22 not trivial [6]. In contrast, humans and other intelligent creatures can learn, without external reward
23 supervision, behaviors that produce repeatable and predictable changes in the environment [17].
24 These behaviors, which we call *skills*, can be later repurposed to solve downstream tasks efficiently.
25 One of the promises of this form of unsupervised RL is to endow artificial agents with similar
26 capabilities to discover reusable skills without explicit rewards.

27 One predominant strategy of prior skill discovery methods focuses on training skills to reach diverse
28 states while being distinguishable [18, 57, 48]. However, in complex environments that contain many
29 *state factors*—distinct elements such as individual objects in a household (a formal description in
30 Sec. 2.1), the exponential number of distinct states makes it impossible to learn skills that cover every
31 state. Consequently, these methods result in simple skills that only change the easy-to-control factors
32 (e.g., in a manipulation task moving the agent itself to diverse positions or manipulating each factor
33 independently), and fail to cover other desirable but challenging behaviors. Unsurprisingly, these
34 simple skills often struggle to solve meaningful tasks, resulting in poor downstream performance.

35 Our key insight is that, given a factored state space, the interactions between state factors can often
36 act as a powerful inductive bias in guiding the learning of useful skills. For example, in a household

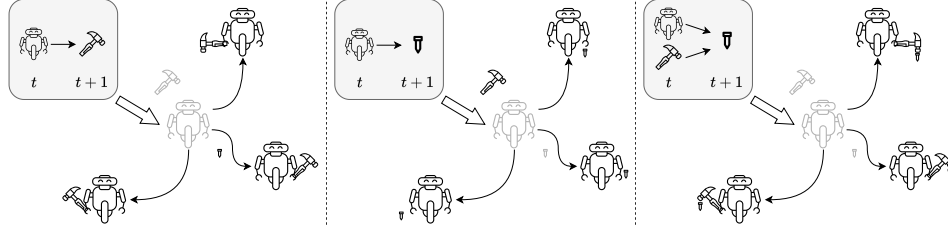


Figure 1: **Skill Discovery from Local Dependencies (SkiLD)** describes skills that encode interactions (i.e., local dependencies) between state factors. In contrast to prior diversity-based methods that can easily get stuck by moving the robot to diverse, but non-interactive states, and factor-based methods that are trained to manipulate the hammer and nail, but not their interactions, SkiLD not only manipulate each object (left, middle) but also induce interactions between them (right), by specifying different local dependencies. These skills are often more useful than the “easy” skill learned by previous methods for downstream task-solving.

37 environment, a skill that induces interactions between a knife and a fruit is more likely to encode
 38 fruit-cutting behaviors that can be crucial for a wide range of downstream kitchen tasks. Furthermore,
 39 exploring state coverage within the space of interaction states can uncover desirable interactions like
 40 cutting the peach.

41 The guiding principle behind this inductive bias is that many domains, including robotic manipulation,
 42 exhibit dynamic bottlenecks through interactions. Interactions act as dynamic bottlenecks by serving
 43 as particular sets of states that must be reached to control another factor. For example, contact between
 44 the knife and fruit is required to control the fruit through cutting. Instead of simply searching for
 45 diversity alone, which in a large state space could focus only on manipulating single factors, forcing
 46 diversity through interactions prevents these dynamic bottlenecks from blocking a wide coverage of
 47 skills.

48 To this end, we introduce Skill Discovery from Local Dependencies (SkiLD), a novel skill discovery
 49 method that explicitly learns skills that induce diverse interactions. Specifically, SkiLD models the
 50 interactions between state factors using the framework of *local dependencies* (where local refers to
 51 state-specific, see details in Sec. 2.2) and proposes a novel intrinsic reward that 1) encourages the agent
 52 to induce specified interactions, and 2) encourages the agent to discover diverse ways of inducing
 53 specified interaction, as visualized in Figure 2. During skill learning, SkiLD gradually discovers
 54 new interactions and learns to induce them, based on the skills that it already mastered, resulting
 55 in a diverse set of interaction-inducing behaviors that can be readily repurposed for downstream
 56 tasks. During task learning, the skill policy is reused, and a task-specific policy is learned to select (a
 57 sequence of) skills to maximize task rewards efficiently.

58 We evaluate the performance of SkiLD on factor-rich environments with 10 downstream tasks against
 59 existing unsupervised reinforcement learning methods. Our experiments indicate that SkiLD learns
 60 to induce diverse interactions and outperforms other methods on most of the examined tasks.

61 2 Background

62 In this paper, our unsupervised skill discovery method is set up in a factored Markov decision process
 63 and builds off previous diversity-based methods, as described in Sec. 2.1. To enhance the expressivity
 64 of skills, our method further augments the skill representation with interactions between state factors,
 65 which we formalize as local dependencies as described in Sec. 2.2.

66 2.1 Factored Markov Decision Process (Factored MDP)

67 We consider unsupervised skill discovery in a reward-free Factored Markov Decision Process [7]
 68 defined by the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p)$. $\mathcal{S} = \mathcal{S}^1 \times \dots \times \mathcal{S}^N$ is a factored state space with N subspaces,
 69 where each subspace \mathcal{S}^i is a multi-dimensional continuous or discrete random variable. Then,
 70 correspondingly, each state $s \in \mathcal{S}$ consists of N state factors, i.e., $s = (s^1, \dots, s^N)$, $s^i \in \mathcal{S}^i$. In this
 71 paper, we use uppercase letters to denote random variables and lowercase for their specific values (e.g.,
 72 S denotes the random variable for states s). \mathcal{A} is the action space, and p is an unknown Markovian
 73 transition model that captures the probability distribution over the next state $S' \sim p(\cdot | S, A)$.

74 The factorization in \mathcal{S} inherently exists in many environments, and is a common assumption in prior
 75 unsupervised skill discovery works [21, 27]. For example, in robotics, an environment typically
 76 consists of a robot and several objects to manipulate, and, for each object, S^i would represent its
 77 attributes of interest, like pose. In this work, we explore how we can utilize a given state factorization
 78 to improve unsupervised skill discovery. In practice, the factorization can either be directly provided
 79 by the environment or obtained from image observations with existing disentangled representation
 80 learning methods [45, 29].

81 Following prior work, our method consists of two stages—skill learning and task learning. During the
 82 skill learning phase, we seek to learn a skill policy $\pi_\omega(\cdot|s, z)$, which defines a conditional distribution
 83 over actions given the current state s and some skill representation z , where skills indicate the desired
 84 behaviors of the agent. Once the skills are learned, they can be chained together to solve downstream
 85 tasks during the task learning phase through an extrinsic reward-optimizing policy. During task
 86 learning, a downstream task reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is provided by the environment. A
 87 high-level policy $\pi(z|s)$ is then trained to optimize the expected return through outputting correct
 88 skills z given state s .

89 2.2 Identifying Local Dependencies between State Factors

90 A key insight of SkiLD is to utilize interactions (or, formally, local dependencies) between state
 91 factors as part of the skill representation. In later sections, these local dependencies are compiled
 92 into a binary matrix $\mathbb{G}(s, a, s') = \{0, 1\}^{N \times (N+1)}$ representing the local dependencies between all
 93 factors. In this section, we first formally define local dependencies, introduce their identification, and
 94 finally discuss their application to factored MDPs.

95 SkiLD takes a causality-inspired approach for defining and detecting local dependencies [5, 56],
 96 where we use *local* to refer to a particular assignment of values for a random variable, as opposed
 97 to *global* which applies to all values. Formally, for an *event of interest* Y and its potential causes
 98 $X = (X^1, \dots, X^N)$, given the value of $X = x$, local dependencies focus on which X^i s are the
 99 state-specific cause of the outcome event $Y = y$ (for simplicity of presentation, in this section we
 100 overload N as the number of potential causes rather than number of variables and p as the transition
 101 function according to a subset of the variables). Formally, we denote the general data generation
 102 process of Y as $p : X \rightarrow Y$ and the data generation process when Y is *only influenced* by a subset of
 103 X as $p^{\bar{X}} : \bar{X} \rightarrow Y$, where $\bar{X} \subseteq X$. Then, given the value of all variables, $X^1 = x^1, \dots, X^N = x^N$
 104 and $Y = y$, we say Y locally depends on \bar{X} , if \bar{X} is the *minimal* subset of X such that knowing their
 105 values is necessary and sufficient to generate the result of $Y = y$, i.e.,

$$\arg \min_{\bar{X} \subseteq X} |\bar{X}| \quad \text{s.t.} \quad p^{\bar{X}}(Y = y | \bar{X} = \bar{x}) = p(Y = y | X = x), \quad (1)$$

106 where $|\bar{X}|$ is the number of variables in \bar{X} . For example, suppose that a robot opens a refrigerator
 107 door in a particular transition. The event of interest Y is the refrigerator door becoming open, and it
 108 locally depends on two factors: the robot and the refrigerator door, while other state factors such as
 109 objects inside the refrigerator do not locally influence Y .

110 To identify local dependencies, one can conduct a conditional independence test $y \perp\!\!\!\perp x^i | \{x/x^i\}$ to
 111 examine whether a variable X^i is necessary for predicting $Y = y$. In prior works, one form of this
 112 test is to examine whether the pointwise conditional mutual information (pCMI) is greater than 0,

$$\text{pCMI}(y; x^i | \{x/x^i\}) = \log \frac{p(y|x)}{p^{\{X/X^i\}}(y | \{x/x^i\})} > 0. \quad (2)$$

113 If so, then it suggests that knowing $X^i = x$ provides additional information about Y that is not
 114 present in $\{X/X^i\}$, and Y locally depends on X^i . As the data generation processes are generally
 115 unknown, one has to approximate them with learned models. Recent work in RL has utilized various
 116 approximations such as attention weights [51], Granger causality [13], and input gradients [60].

117 In this work, for a transition ($\mathcal{S} = s, \mathcal{A} = a, \mathcal{S}' = s'$), the event of interest is each next state factor
 118 being $(\mathcal{S}^i)' = (s^i)'$, and we infer whether it locally depends on each state factor \mathcal{S}^j and the action \mathcal{A}
 119 (i.e., whether there is an interaction between state factors i and j , where factor j influences i). Then
 120 we aggregate all local dependencies into a state-specific dependency graph (abbreviated in this work to
 121 *dependency graph*). This overall dependency graph is represented with $\mathbb{G}(s, a, s') = \{0, 1\}^{N \times (N+1)}$,

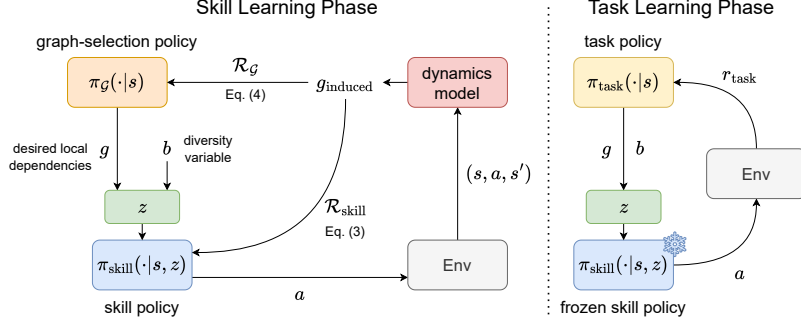


Figure 2: During **skill learning** of SkiLD, the graph-selection policy specifies desired local dependencies for the skill policy to induce, and the induced dependency graph is identified by the dynamics model and used to update both policies. During **task learning** (right), the skill policy is kept frozen and a task policy is trained to select skills to maximize task reward.

122 and an edge $\mathbb{G}^{ij}(s, a, s')$ denotes, during the transition (s, a, s') , that state factor $(s^i)'$ (the “ $Y = y$ ”)
 123 locally depends on s^j (one of the X^j):

$$\mathbb{G}^{ij} := \text{pCMI}((x^i)'; x^j | \{x/x^j\}) \quad (3)$$

124 This graph is used to enhance skill representation, as explained in detail in Section 3.

125 3 Skill Discovery from Local Dependencies (SkiLD)

126 In this section, we describe SkiLD, which enhances the expressivity of skills using local dependencies.
 127 SkiLD represents local dependencies as *state-specific dependency graphs*, defined in Sec. 2.2. Unlike
 128 previous unsupervised skill discovery methods that randomly sample the skill vector z from fixed
 129 distributions during skill learning, SkiLD requires a procedure to intelligently generate target dependency
 130 graphs during training. As such, SkiLD frames unsupervised skill discovery as a hierarchical
 131 RL problem, where a graph-conditioned skill policy learns to induce different local dependencies
 132 using primitive actions, and a high-level graph selection policy chooses which local dependencies the
 133 skill policy should induce next to guide exploration and skill-policy learning.

134 This requires formalizing two components: (1) the skill representation \mathcal{Z} for the skill policy
 135 $\pi_{\text{skill}}(a|s, z)$ and its corresponding reward function $\mathcal{R}_{\text{skill}}$, presented in Sec. 3.1, and (2) the graph
 136 selection policy $\pi_{\mathbb{G}}(z|s)$ and its reward function $\mathcal{R}_{\mathbb{G}}$, presented in Sec. 3.2.

137 3.1 Skill Policy

138 Prior unsupervised skill discovery methods usually focus skill learning on changing the state or each
 139 factor diversely. Consequently, they can be limited to learning simple skills, for example, only
 140 changing the easiest-to-control factor in the state (i.e., the agent itself). To address this problem,
 141 SkiLD not only focuses on changing the state but also considers the interactions between state factors.

142 **Skill Representation.** SkiLD represents the skill space as the combination of two components:
 143 $\mathcal{Z} = \mathbb{G} \times \mathcal{B}$, where $g \in \mathcal{G}$ is a state-specific dependency graph that specifies the *desired* local
 144 dependencies between state factors, and $b \in \mathcal{B}$ is a diversity variable the same as that used in
 145 Eysenbach et al. [18]. Together $z \in \mathcal{Z}$ guides the agent to change the state distinguishably while
 146 inducing particular local dependencies. Specifically, the dependency graph is represented as a
 147 binary matrix $\mathbb{G} = \{0, 1\}^{N \times (N+1)}$, where each edge \mathbb{G}^{ij} denotes, during the transition (s, a, s') ,
 148 whether the state factor $(s^i)'$ locally depends on s^j . The diversity variable \mathcal{B} can be either discrete or
 149 continuous. In this work, without loss of generality, we use a discrete space of $\{1, \dots, K\}$ where K
 150 is a predefined number. During skill training, we sample the diversity variable b from a fixed uniform
 151 distribution $p(b)$, following the procedure of Eysenbach et al. [18].

152 Given this skill space, SkiLD learns skills as a skill-conditioned policy $\pi_{\text{skill}} : \mathcal{S} \times \mathcal{Z} \rightarrow \mathcal{A}$, where
 153 π_{skill} is trained to reach diverse states while ensuring that the local dependencies specified by the
 154 graph are induced. During skill learning, we select actions by iteratively calling the skill policy
 155 π_{skill} , and we denote g_{induced} as the graph that describes the local dependencies induced in a transition

156 (s, a, s') when executing a selected action a . We design the reward function of SkiLD as:

$$\mathcal{R}_{\text{skill}} = \mathbb{1}[g_{\text{induced}} = g] \cdot (1 + \lambda \mathcal{R}_{\text{diversity}}), \quad (4)$$

157 where $\mathbb{1}[g_{\text{induced}} = g]$ measures whether the induced dependency graph matches the desired graph,
 158 $\mathcal{R}_{\text{diversity}}$ is the weighted diversity reward that further encourages visiting diverse states when the
 159 desired graph is induced, and λ is the coefficient of diversity reward. In the following paragraphs, we
 160 describe how we infer g_{induced} and estimate $\mathcal{R}_{\text{diversity}}$ for each transition.

161 **Inferring Induced Graphs.** To infer the induced graph for a transition $(S = s, A = a, S' = s')$,
 162 we need to determine, for each $(S')^i$, whether it locally depends on each factor S^j and the action
 163 \mathcal{A} . Specifically, following Sec. 2.2, we evaluate the conditional dependency $(s^i)' \not\perp\!\!\!\perp s^j | \{s/s^j, a\}$
 164 by examining whether their pointwise conditional mutual information (pCMI) is greater than a
 165 predefined threshold $\text{pCMI}^{i,j} = \frac{p((s^i)'|s,a)}{p((s^i)'|\{s/s^j,a\})} \geq \epsilon$. If so, it suggests that s^j is necessary to predict
 166 $(s^i)'$ and thus the local dependency exists. Meanwhile, as the transition probability p is unknown, we
 167 approximate it with a learned dynamics model that is trained to minimize prediction error.

168 Finally, after obtaining the induced dependency graph, we evaluate $\mathbb{1}[g_{\text{induced}} = g]$ by examining
 169 whether each edge $g_{\text{induced}}^{i,j}$ matches the corresponding edge in the desired graph $g^{i,j}$. As $\mathcal{R}_{\text{skill}}$ only
 170 provides sparse rewards to the skill policy when the desired graph is induced, we use hindsight
 171 experience replay [1] to enrich learning signals, by relabelling induced graphs as desired graphs in
 172 some episodes.

173 **Diversity Rewards.** When the skill policy induces the desired graph, $\mathcal{R}_{\text{diversity}}$ further encourages
 174 it to visit different distinguishable states under different diversity indicators b , e.g., driving the nail
 175 to different locations. This diversity enhances the applicability of learned skills. To this end, we
 176 design the diversity reward $\mathcal{R}_{\text{diversity}}$ as the forward mutual information between visited states and the
 177 diversity indicator $I(s; b)$, following DIAYN. To estimate the mutual information, we approximate
 178 it with a variational lower bound $I(s; b) \geq q(b|s)$, where $q(b|s)$ is a neural network discriminator
 179 trained to predict the diversity indicator b from the visited state. In practice, rather than learning a
 180 single low-level skill to handle all graphs and diversity parameters, we utilize a factorized lower-level
 181 policy, where there is a separate policy for each factor. More details about this subdivision can be
 182 found in Appendix A.

183 3.2 Graph-Selection Policy

184 To acquire skills that are useful for downstream tasks, π_{skill} needs to learn to induce a wide range
 185 of local dependencies *sample-efficiently*. To this end, we propose to learn a graph-selection policy
 186 $\pi_{\mathbb{G}} : \mathcal{S} \rightarrow \mathbb{G}$ to guide the training of π_{skill} . Specifically, training π_{skill} requires a wise selection of
 187 graphs — as graph space \mathbb{G} increases super-exponentially in the number of state factors N , many
 188 graphs are not inducible. To this end, we only select target graphs for skill policy from a history of
 189 all seen graphs. As the agent learns to induce existing graphs in diverse ways, new graphs may be
 190 encountered, gradually expanding the set of seen graphs.

191 However, though this history guarantees graph inducibility, two challenges still remain: (1) How to
 192 efficiently explore novel local dependencies, especially hard-to-visit ones? (2) For all seen graphs,
 193 which one should π_{skill} learn next to maximize training efficiency? We address these challenges
 194 based on the following heuristic — compared to well-learned skills, π_{skill} should focus its training on
 195 underdeveloped skills. Meanwhile, learning new skills opens up the possibility of visiting novel local
 196 dependencies, e.g., learning to grasp the hammer makes it possible for the robot to hammer the nail.

197 According to this heuristic, we learn a graph-selection policy $\pi_{\mathbb{G}}$ that guides the exploration and
 198 training of the skill policy π_{skill} . Specifically, $\pi_{\mathbb{G}} : \mathcal{S} \rightarrow \mathbb{G}$ selects a new dependency graph the skill
 199 policy should induce for the next L time steps. To increase the likelihood of visiting hard graphs, $\pi_{\mathbb{G}}$
 200 is trained to maximize the following graph novelty reward

$$\mathcal{R}_{\mathbb{G}} = \frac{1}{\sqrt{C(g_{\text{visited}})}}, \quad (5)$$

201 where $C(g_{\text{visited}})$ is the number of times that we have seen the graph in the collected transition. While
 202 Eq. 5 is similar to state-count-based exploration reward, here, it is based on the count of dependency
 203 graphs, and thus applicable to both discrete and continuous state space.

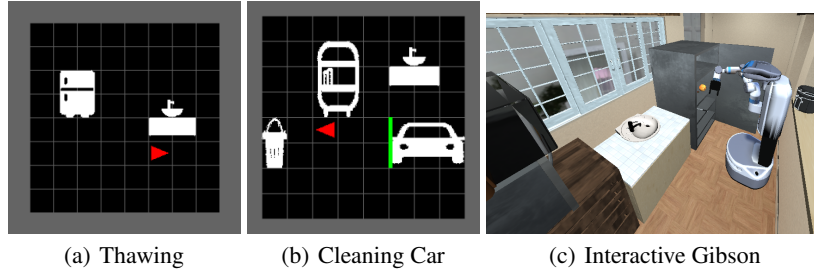


Figure 3: **Evaluation domains:** Mini-behavior: Installing Printer, Thawing and Cleaning Car, and iGibson.

204 3.3 Downstream Task Learning

205 In SkiLD, we utilize hierarchical RL to solve reward-supervised downstream tasks with the discovered
 206 skills. The skill policy, π_{skill} acts as the low-level policy while a task policy, $\pi_{\text{task}} : \mathcal{S} \rightarrow \mathcal{Z}$, is learned
 207 to select which skill $z = (g, b)$ to execute for L steps. Compared to diversity-based skills that are
 208 limited to simple behaviors, our local-dependency-based skills enable a wide range of interactions
 209 between state factors, leading to more efficient exploration and superior performance of downstream
 210 tasks learning.

211 4 Experiments

212 In this section we aim to provide empirical evidence towards the following questions: **Q1**) Do
 213 the skills learned by SkiLD induce a diverse set of interactions among state factors? **Q2**) Do
 214 the skills learned by SkiLD enable more efficient downstream task learning compared to other
 215 unsupervised reinforcement learning methods? Our learned skills can be visualized at <https://sites.google.com/view/skild/>.
 216

217 4.1 Domains

218 In this work, we focus on addressing the challenge of vast state space brought by a large number
 219 of state factors. Hence, we evaluate our method on two challenging *object-rich* embodied AI
 220 benchmarks: Mini-behavior [30] and Interactive Gibson [40].

221 The **Mini-behavior (Mini-BH) domain** [30] (Figure 3a) contains a set of gridworld environments
 222 where an agent can move around and interact with a variety of objects to accomplish certain household
 223 tasks. While conceptually simple, this domain has been shown to be extremely challenging for Vanilla
 224 RL with sparse reward [30]. Each Mini-BH environment contains different objects and different
 225 success criteria. We tested on three particular environments in Mini-behavior, including:

- 226 • **Installing Printer:** A relatively simple environment with three state factors: the agent, a table, and
 227 a printer that can be installed.
- 228 • **Cleaning Car:** An environment where the objects have rich and complex interactions. The state
 229 factors include the agent, a toggleable sink, a piece of rag that can be soaked in the sink, a car that
 230 the rag can clean, a soap and a bucket which can together be used to clean the rag.
- 231 • **Thawing:** An environment with lots of movable objects. The state factors include the agent, a sink,
 232 a fridge that can be opened, and three objects that can be thawed in the sink: fish, olive, and a date.

233 The **Interactive Gibson (iGibson) domain** [41] (Figure 3b) contains a realistic simulated Fetch
 234 Robot that operates in a kitchen environment with a refrigerator, sink, knife, and peach. The peach
 235 can be washed or cut. This domain is very difficult especially when using low-level motor commands
 236 because much of the domain is free space, meaning that only a minute fraction of action sequences
 237 will manipulate the objects meaningfully.

238 Both Mini-BH and iGibson require learning long-horizon policies spanning many low-level actions
 239 from sparse reward, making these challenging domains (see details in Appendix).

240 **4.2 Baselines**

241 Before evaluating the empirical questions, we provide a brief description of the baselines. These
 242 baselines include unsupervised skill learning, and causal and hierarchical methods.

243 **Diversity is all you need (DIAYN [18]):** This method learns unsupervised state-covering skills using
 244 a mutual information objective. SkiLD utilizes a version of this for state-diversity skills modulated by
 245 a desired dependency graph. This baseline determines how incorporating graph information affects
 246 the algorithm.

247 **Controllability-Aware Skill Discovery (CSD [48]):** Extends DIAYN with a factorization based on
 248 controllability. This baseline is a comparable skill learning method that leverages state factorization
 249 but does not encode local dependencies.

250 **Exploration via Local Dependencies (ELDEN [60]):** This method utilizes gradient-based techniques
 251 to infer local dependencies for exploration. However, without a skill learning component, it can
 252 struggle to chain together complex behavior.

253 **Chain of Interaction Skills (COInS [13]):** This is a hierarchical algorithm that constructs a chain
 254 of skills using Granger-causality to identify local dependencies. Because it is restricted to pairwise
 255 interactions, it struggles to represent the rich policies necessary for these tasks.

256 **Vanilla RL:** This baseline uses PPO [55] to directly train an agent with the extrinsic reward. Unlike
 257 other baselines, this method does not have a pertaining phase. Since all the task rewards are sparse
 258 and the tasks are often long horizon, vanilla RL often struggles.

259 **4.3 Interaction Graph Diversity**

260 We first evaluate whether SkiLD is indeed ca-
 261 pable of achieving complex interaction graphs
 262 (Q1), comparing against two strong skill discov-
 263 ery baselines introduced earlier: DIAYN and
 264 CSD.

265 Each of these methods is trained for 10 Million
 266 steps without having access to any reward. Then
 267 to evaluate their learned skills, we unroll each
 268 of them for 500 episodes with randomly sam-
 269 pled skills z and examine the diversity of the
 270 interaction graphs they can induce. Figure 4
 271 illustrates the percentages of episodes where
 272 particular local dependencies have been induced
 273 at least once, in Mini-BH Cleaning Car. We
 274 find that DIAYN and CSD are limited to skills
 275 that only manipulate one object individually, i.e.
 276 (agent, rag, action \rightarrow rag) or (agent, soap,
 277 action \rightarrow soap). By contrast, SkiLD learns to
 278 induce more complicated causal interactions, such
 279 as soaking the rag in the sink (sink, rag \rightarrow rag)
 280 and cleaning the car with the soaked mug (car,
 281 rag \rightarrow car).

282 **4.4 Performance**

283 Next, we evaluate whether the local dependency coverage provided by SkiLD leads to a performance
 284 boost in downstream task learning (Q2). We follow the evaluation setup in the unsupervised rein-
 285 forcement learning benchmark [36], where for a given environment, an agent is first pre-trained
 286 without access to task reward for K_{pt} steps, and then finetuned for K_{ft} steps. Importantly, the same
 287 pre-trained skills are reused on multiple distinct downstream tasks within the same environment, so
 288 that only the upper-level skill-selection policy is task-specific. We have $K_{pt} = 2M$, $K_{ft} = 1M$ for
 289 installing printer, $K_{pt} = 10M$, $K_{ft} = 5M$ for thawing and cleaning car, and $K_{pt} = 4M$, $K_{ft} = 2M$

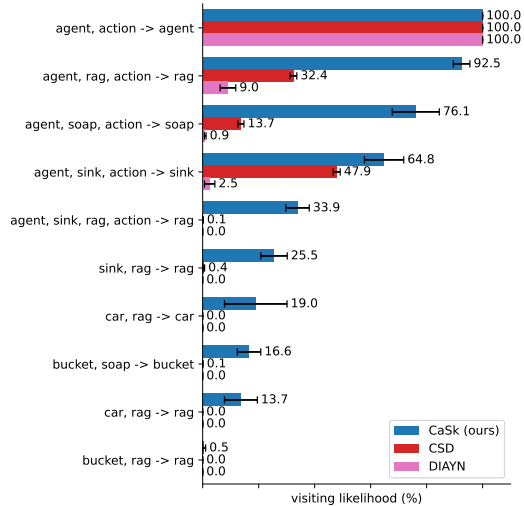


Figure 4: The percentage of episodes where a dependency graph is induced through random skill sampling. Standard deviation is calculated across five random seeds.

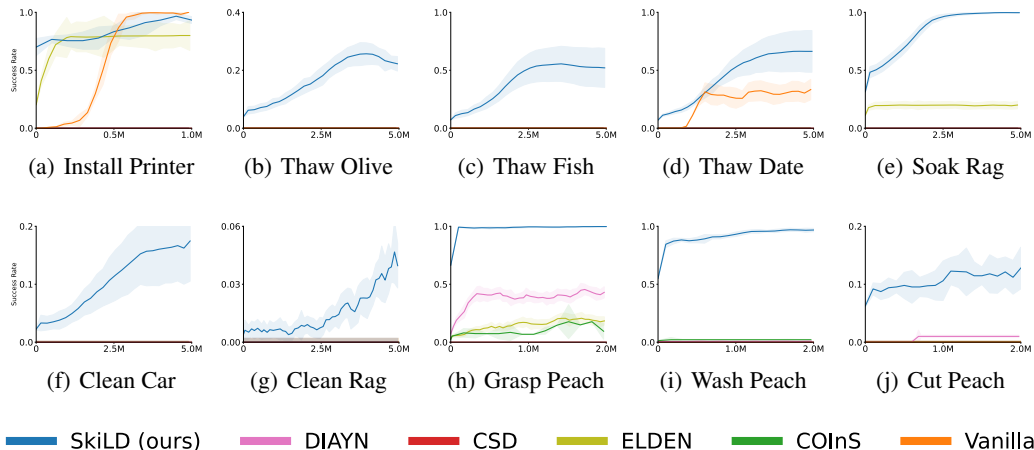


Figure 5: Training curves of SkiLD and baselines on multiple downstream tasks (reward supervised second phase). Each curve depicts the mean and standard deviation of the success rate over 5 random seeds. SkiLD outperforms all baselines for most tasks, converging faster and to higher returns.

290 for iGibson, and evaluate each method for each task across 5 random seeds. Hyperparameter details
 291 can be found in Appendix D. Specifically, we evaluate on the following downstream tasks:

- 292 • **Installing Printer:** We have a single downstream task in this environment, where the agent needs
 293 to pick up the printer, put it on the table, and turn it on.
- 294 • **Thawing:** We have three downstream tasks: thawing the fish or the olive or the date.
- 295 • **Cleaning Car:** We consider three downstream tasks, where each task is a pre-requisite of the
 296 following one. The tasks are: soak the rag in the sink; clean the car with the rag; and clean the dirty
 297 rag using the soap in the bucket.
- 298 • **iGibson:** The tasks for this domain are: grasping the peach, washing the peach in the sink, and
 299 cutting the peach with a knife.

300 After skill learning, we train a new upper-level policy that uses z as actions and is trained with extrinsic
 301 reward, as described in Section 3.3. Figure 5 illustrates the improvement of SkiLD as compared to
 302 other methods. Without combining dependency graphs with skill learning, other methods struggle
 303 with any but the simpler tasks. COInS performs poorly because of its chain structure, which restricts
 304 the agent controlling policy from picking up objects. ELDEN’s exploration reaches graphs, but
 305 without skills struggles to utilize that information in downstream tasks. DIAYN learns skills, but few
 306 manipulate the objects, so a downstream model struggles to utilize those skills to achieve meaningful
 307 rewards. By comparison, SkiLD achieves superior performance on 9 of the 10 downstream tasks
 308 evaluated. In the two hardest tasks which require a very long sequence of precise controls, Clean Rag
 309 and Cut Peach, SkiLD is the only method that can achieve a non-zero success rate (although still far
 310 from fully mastering the tasks), showcasing the potential of local dependencies for skill learning.

311 4.5 Graph and Diversity Ablations

312 We also explore the functionality of the graph and diversity components of the skill parameter z
 313 by assessing the downstream performance of SkiLD without these components. This produces two
 314 ablative versions of SkiLD: SkiLD without diversity and SkiLD without dependency graphs. To
 315 isolate learning from the effect of learned local dependencies, we use ground truth dependency
 316 graphs for ablative evaluations where relevant. In Figure 6, learning without graphs results in zero
 317 performance, consistent with DIAYN results. In addition, removing diversity produces a notable
 318 decline in performance, especially on more challenging tasks like cleaning the rag. These evaluations
 319 demonstrate that SkiLD benefits from both the incorporation of dependency graphs and diversity.

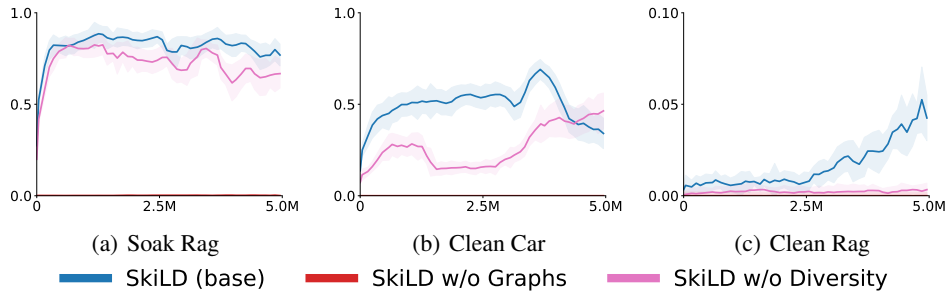


Figure 6: A figure illustrating the ablative performance of SkiLD without diversity or without graphs. Without graphs, the method collapses completely, while removing diversity results in a noticeable reduction in downstream performance.

320 5 Related Work

321 This work lies in the unsupervised skill learning framework [35], where the agent must discover a
 322 set of useful skills which are reward independent. It then extends these skills to construct a 2-layer
 323 hierarchical structure [58], where the upper policy receives reward both for achieving novel skills,
 324 and can then be tuned to utilize the learned skills to accomplish an end task. Finally, the skills are
 325 identified using token causality, a specific problem identified in causal literature.

326 5.1 Unsupervised Skill Learning

327 This work describes a framework for utilizing local dependency graphs and diversity to discover
 328 unsupervised skills. Diversity-based state coverage skills have been explored in literature [18]
 329 utilizing forward and backward mutual information techniques to learn a goal space \mathcal{Z} , and a skill
 330 encoder $q(z|\cdot)$ [10]. This unsupervised paradigm has been extended with Lipschitz constraints [47],
 331 contrastive objectives [37], information bottleneck [33], population based methods such as particle
 332 estimation [43], quality diversity [42] and mixture of experts [11]. These skills can then be used for
 333 hierarchical policies or planners [54, 64, 22], which mirrors the same structure as SkiLD. Unlike
 334 these methods, SkiLD adds additional subdivision through dependency graphs, which mitigates the
 335 combinatorial explosion of skills that can result from trying to cover a large factored space.

336 5.2 Hierarchical Reinforcement Learning

337 The hierarchical policy structure in SkiLD where a higher level policy passes a parameter to be inter-
 338 preted by low-level planners has been formalized in [58], and learned using deep networks utilizing
 339 extrinsic reward [2, 59], attention mechanisms [15], initiation criteria [32, 3] and deliberation cost [25].
 340 Hierarchies of goal-based policies [38] has been extended with object-centric representations [63],
 341 offline data [46], empowerment [39] and goal counts [49]. In practice, SkiLD uses graph and diversity
 342 parameters similar to goal-based methods. However, the space of goals can often be intractable
 343 large, and methods to address this use graph laplacians [34] causal chains [12, 13] or general causal
 344 relationships [27]. SkiLD is similar to these causal methods but utilizes local dependence along with
 345 general two-layer architectures, thus showing increased generalizability.

346 5.3 Causality in Reinforcement Learning

347 This work investigates the application of local dependency to hierarchical reinforcement learning.
 348 This kind of reasoning has been described as “local causality” or “interactions” in prior RL work
 349 for data augmentation [51, 52], learning skill chains [12, 13] and exploration [60]. This work is
 350 the first synthesis of unsupervised skill learning and local dependencies applied to general 2-layer
 351 hierarchical reinforcement learning. Other general causality work investigates action-influence
 352 detection [56, 26], affordance learning [9], model learning [28, 20], critical state identification [44],
 353 and disentanglement [16]. In the context of relating local dependency and causal inference, we
 354 provide a discussion in Appendix C. SkiLD incorporates causality-inspired local dependence to skill
 355 learning, resulting in a robust set of transferable skills.

356 6 Conclusion

357 Unsupervised skill discovery is a powerful tool for learning useful skills in long-horizon sparse
358 reward tasks. However, many unsupervised skill-learning methods do not take advantage of factored
359 environments, resulting in poor performance in complex environments with several objects. Skill
360 Discovery from Local Dependencies utilizes state-specific dependency graphs, identified using
361 learned pointwise conditional mutual information models, to guide skill discovery. The framework
362 of defining skills according to a dependency graph and diversity goal, combined with a learned
363 sampling scheme, achieves difficult downstream tasks. In domains where hand-coded primitive
364 skills are typically given to the agent, like Mini-behavior and Interactive Gibson, SkiLD can achieve
365 high performance without requiring explicit domain knowledge. These impressive results arise
366 intuitively from incorporating local dependencies as skill targets, illuminating a meaningful direction
367 for unsupervised skill learning to be applied to a wider array of environments.

368 **Limitations and Future Work** An important assumption of SkiLD is its access to factored state
369 space. While factored state space can often be naturally obtained from existing RL benchmarks
370 and many real-world environments, developments in disentangled representation learning [45, 29]
371 will help with extending SkiLD to unfactored image domains. Secondly, SkiLD requires accurate
372 detection of local dependencies. While off-the-shelf methods [60, 56] work well for detecting local
373 dependencies in our experiments, future works that can more accurately detect local dependencies
374 will be beneficial to the performance of SkiLD.

375 References

- 376 [1] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder,
377 Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience
378 replay. *Advances in neural information processing systems*, 30, 2017.
- 379 [2] Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Proceedings*
380 *of the AAAI conference on artificial intelligence*, volume 31, 2017.
- 381 [3] Akhil Bagaria, Ben Abbatematteo, Omer Gottesman, Matt Corsaro, Sreehari Rammohan, and
382 George Konidaris. Effectively learning initiation sets in hierarchical reinforcement learning.
383 *Advances in Neural Information Processing Systems*, 36, 2024.
- 384 [4] Sander Beckers, Hana Chockler, and Joseph Halpern. A causal analysis of harm. *Advances in*
385 *Neural Information Processing Systems*, 35:2365–2376, 2022.
- 386 [5] Gianluca Bontempi and Maxime Flauder. From dependency to causality: a machine learning
387 approach. *J. Mach. Learn. Res.*, 16(1):2437–2457, 2015.
- 388 [6] Serena Booth, W Bradley Knox, Julie Shah, Scott Niekum, Peter Stone, and Alessandro Allievi.
389 The perils of trial-and-error reward design: misdesign through overfitting and invalid task
390 specifications. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37,
391 pages 5920–5929, 2023.
- 392 [7] Craig Boutilier, Thomas Dean, and Steve Hanks. Decision-theoretic planning: Structural
393 assumptions and computational leverage. *Journal of Artificial Intelligence Research*, 11:1–94,
394 1999.
- 395 [8] Craig Boutilier, Nir Friedman, Moises Goldszmidt, and Daphne Koller. Context-specific
396 independence in bayesian networks. *arXiv preprint arXiv:1302.3562*, 2013.
- 397 [9] Jake Brawer, Meiyang Qin, and Brian Scassellati. A causal approach to tool affordance learning.
398 In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages
399 8394–8399. IEEE, 2020.
- 400 [10] Víctor Campos, Alexander Trott, Caiming Xiong, Richard Socher, Xavier Giró-i Nieto, and
401 Jordi Torres. Explore, discover and learn: Unsupervised discovery of state-covering skills. In
402 *International Conference on Machine Learning*, pages 1317–1327. PMLR, 2020.

- 403 [11] Onur Celik, Dongzhuoran Zhou, Ge Li, Philipp Becker, and Gerhard Neumann. Specializing
404 versatile skill libraries using local mixture of experts. In *Conference on Robot Learning*, pages
405 1423–1433. PMLR, 2022.
- 406 [12] Caleb Chuck, Supawit Chockchawat, and Scott Niekum. Hypothesis-driven skill discovery
407 for hierarchical deep reinforcement learning. In *2020 IEEE/RSJ International Conference on*
408 *Intelligent Robots and Systems (IROS)*, pages 5572–5579. IEEE, 2020.
- 409 [13] Caleb Chuck, Kevin Black, Aditya Arjun, Yuke Zhu, and Scott Niekum. Granger-causal
410 hierarchical skill discovery. *arXiv preprint arXiv:2306.09509*, 2023.
- 411 [14] Caleb Chuck, Sankaran Vaidyanathan, Stephen Giguere, Amy Zhang, David Jensen, and Scott
412 Niekum. Automated discovery of functional actual causes in complex environments. *arXiv*
413 *preprint arXiv:2404.10883*, 2024.
- 414 [15] Raviteja Chunduru and Doina Precup. Attention option-critic. *arXiv preprint arXiv:2201.02628*,
415 2022.
- 416 [16] Oriol Corcoll and Raul Vicente. Disentangling controlled effects for hierarchical reinforcement
417 learning. In Bernhard Schölkopf, Caroline Uhler, and Kun Zhang, editors, *Proceedings of the*
418 *First Conference on Causal Learning and Reasoning*, volume 177 of *Proceedings of Machine*
419 *Learning Research*, pages 178–200. PMLR, 11–13 Apr 2022. URL [https://proceedings.](https://proceedings.mlr.press/v177/corcoll122a.html)
420 [mlr.press/v177/corcoll122a.html](https://proceedings.mlr.press/v177/corcoll122a.html).
- 421 [17] Yuqing Du, Eliza Kosoy, Alyssa Dayan, Maria Rufova, Pieter Abbeel, and Alison Gopnik. What
422 can ai learn from human exploration? intrinsically-motivated humans and agents in open-world
423 exploration. In *NeurIPS 2023 workshop: Information-Theoretic Principles in Cognitive Systems*,
424 2023.
- 425 [18] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you
426 need: Learning skills without a reward function. *arXiv preprint arXiv:1802.06070*, 2018.
- 427 [19] Alhussein Fawzi, Matej Balog, Aja Huang, Thomas Hubert, Bernardino Romera-Paredes,
428 Mohammadamin Barekatin, Alexander Novikov, Francisco J R Ruiz, Julian Schrittwieser,
429 Grzegorz Swirszcz, et al. Discovering faster matrix multiplication algorithms with reinforcement
430 learning. *Nature*, 610(7930):47–53, 2022.
- 431 [20] Fan Feng and Sara Magliacane. Learning dynamic attribute-factored world models for efficient
432 multi-object reinforcement learning. *Advances in Neural Information Processing Systems*, 36,
433 2024.
- 434 [21] Pierre Fournier, Cédric Colas, Mohamed Chetouani, and Olivier Sigaud. Clic: Curriculum
435 learning and imitation for object control in nonrewarding environments. *IEEE Transactions on*
436 *Cognitive and Developmental Systems*, 13(2):239–248, 2019.
- 437 [22] Tuomas Haarnoja, Kristian Hartikainen, Pieter Abbeel, and Sergey Levine. Latent space policies
438 for hierarchical reinforcement learning. In *International Conference on Machine Learning*,
439 pages 1851–1860. PMLR, 2018.
- 440 [23] Joseph Y Halpern. *Actual causality*. MIT Press, 2016.
- 441 [24] Joseph Y Halpern and Judea Pearl. Causes and explanations: A structural-model approach. part
442 i: Causes. *The British journal for the philosophy of science*, 2005.
- 443 [25] Jean Harb, Pierre-Luc Bacon, Martin Klissarov, and Doina Precup. When waiting is not an
444 option: Learning options with a deliberation cost. In *Proceedings of the AAAI Conference on*
445 *Artificial Intelligence*, volume 32, 2018.
- 446 [26] Jiaheng Hu, Peter Stone, and Roberto Martín-Martín. Causal policy gradient for whole-body
447 mobile manipulation. *arXiv preprint arXiv:2305.04866*, 2023.
- 448 [27] Xing Hu, Rui Zhang, Ke Tang, Jiaming Guo, Qi Yi, Ruizhi Chen, Zidong Du, Ling Li, Qi Guo,
449 Yunji Chen, et al. Causality-driven hierarchical structure discovery for reinforcement learning.
450 *Advances in Neural Information Processing Systems*, 35:20064–20076, 2022.

- 451 [28] Yixuan Huang, Adam Conkey, and Tucker Hermans. Planning for multi-object manipulation
452 with graph neural network relational classifiers. In *2023 IEEE International Conference on*
453 *Robotics and Automation (ICRA)*, pages 1822–1829. IEEE, 2023.
- 454 [29] Jindong Jiang, Fei Deng, Gautam Singh, and Sungjin Ahn. Object-centric slot diffusion. *arXiv*
455 *preprint arXiv:2303.10834*, 2023.
- 456 [30] Emily Jin, Jiaheng Hu, Zhuoyi Huang, Ruohan Zhang, Jiajun Wu, Li Fei-Fei, and Roberto
457 Martín-Martín. Mini-behavior: A procedurally generated benchmark for long-horizon decision-
458 making in embodied ai. *arXiv preprint arXiv:2310.01824*, 2023.
- 459 [31] Elia Kaufmann, Leonard Bauersfeld, Antonio Loquercio, Matthias Müller, Vladlen Koltun, and
460 Davide Scaramuzza. Champion-level drone racing using deep reinforcement learning. *Nature*,
461 620(7976):982–987, 2023.
- 462 [32] Khimya Khetarpal, Martin Klissarov, Maxime Chevalier-Boisvert, Pierre-Luc Bacon, and Doina
463 Precup. Options of interest: Temporal abstraction with interest functions. In *Proceedings of the*
464 *AAAI Conference on Artificial Intelligence*, volume 34, pages 4444–4451, 2020.
- 465 [33] Jaekyeom Kim, Seohong Park, and Gunhee Kim. Unsupervised skill discovery with bottleneck
466 option learning. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International*
467 *Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*,
468 pages 5572–5582. PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.press/v139/
469 kim21j.html](https://proceedings.mlr.press/v139/kim21j.html).
- 470 [34] Martin Klissarov and Marlos C. Machado. Deep Laplacian-based options for temporally-
471 extended exploration. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt,
472 Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International*
473 *Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*,
474 pages 17198–17217. PMLR, 23–29 Jul 2023. URL [https://proceedings.mlr.press/
475 v202/klissarov23a.html](https://proceedings.mlr.press/v202/klissarov23a.html).
- 476 [35] Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. Exploration in deep re-
477 inforcement learning: A survey. *Information Fusion*, 85:1–22, 2022. ISSN 1566-2535.
478 doi: <https://doi.org/10.1016/j.inffus.2022.03.003>. URL [https://www.sciencedirect.com/
479 science/article/pii/S1566253522000288](https://www.sciencedirect.com/science/article/pii/S1566253522000288).
- 480 [36] Michael Laskin, Denis Yarats, Hao Liu, Kimin Lee, Albert Zhan, Kevin Lu, Catherine Cang,
481 Lerrel Pinto, and Pieter Abbeel. Urlb: Unsupervised reinforcement learning benchmark, 2021.
- 482 [37] Michael Laskin, Hao Liu, Xue Bin Peng, Denis Yarats, Aravind Rajeswaran, and Pieter
483 Abbeel. Cic: Contrastive intrinsic control for unsupervised skill discovery. *arXiv preprint*
484 *arXiv:2202.00161*, 2022.
- 485 [38] Andrew Levy, Robert Platt, and Kate Saenko. Hierarchical reinforcement learning with hindsight.
486 In *International Conference on Learning Representations*, 2019. URL [https://openreview.
487 net/forum?id=ryzECoAcY7](https://openreview.net/forum?id=ryzECoAcY7).
- 488 [39] Andrew Levy, Sreehari Rammohan, Alessandro Allievi, Scott Niekum, and George Konidaris.
489 Hierarchical empowerment: Towards tractable empowerment-based skill-learning. *arXiv*
490 *preprint arXiv:2307.02728*, 2023.
- 491 [40] Chengshu Li, Fei Xia, Roberto Martín-Martín, Michael Lingelbach, Sanjana Srivastava, Bokui
492 Shen, Kent Vainio, Cem Gokmen, Gokul Dharan, Tanish Jain, Andrey Kurenkov, C. Karen
493 Liu, Hyowon Gweon, Jiajun Wu, Li Fei-Fei, and Silvio Savarese. igibson 2.0: Object-centric
494 simulation for robot learning of everyday household tasks, 2021.
- 495 [41] Chengshu Li, Fei Xia, Roberto Martín-Martín, Michael Lingelbach, Sanjana Srivastava, Bokui
496 Shen, Kent Elliott Vainio, Cem Gokmen, Gokul Dharan, Tanish Jain, Andrey Kurenkov, Karen
497 Liu, Hyowon Gweon, Jiajun Wu, Li Fei-Fei, and Silvio Savarese. igibson 2.0: Object-centric
498 simulation for robot learning of everyday household tasks. In Aleksandra Faust, David Hsu,
499 and Gerhard Neumann, editors, *Proceedings of the 5th Conference on Robot Learning*, volume
500 164 of *Proceedings of Machine Learning Research*, pages 455–465. PMLR, 08–11 Nov 2022.
501 URL <https://proceedings.mlr.press/v164/li22b.html>.

- 502 [42] Bryan Lim, Luca Grillotti, Lorenzo Bernasconi, and Antoine Cully. Dynamics-aware quality-
503 diversity for efficient learning of skill repertoires. In *2022 International Conference on Robotics*
504 *and Automation (ICRA)*, pages 5360–5366. IEEE, 2022.
- 505 [43] Hao Liu and Pieter Abbeel. Behavior from the void: Unsupervised active pre-training. *Advances*
506 *in Neural Information Processing Systems*, 34:18459–18473, 2021.
- 507 [44] Haozhe Liu, Mingchen Zhuge, Bing Li, Yuhui Wang, Francesco Faccio, Bernard Ghanem,
508 and Jürgen Schmidhuber. Learning to identify critical states for reinforcement learning from
509 videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages
510 1955–1965, 2023.
- 511 [45] Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and
512 Michael Tschannen. Weakly-supervised disentanglement without compromises. In *International*
513 *Conference on Machine Learning*, pages 6348–6359. PMLR, 2020.
- 514 [46] Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical
515 reinforcement learning. *Advances in neural information processing systems*, 31, 2018.
- 516 [47] Seohong Park, Jongwook Choi, Jaekyeom Kim, Honglak Lee, and Gunhee Kim. Lipschitz-
517 constrained unsupervised skill discovery. In *International Conference on Learning Representa-*
518 *tions*, 2021.
- 519 [48] Seohong Park, Kimin Lee, Youngwoon Lee, and Pieter Abbeel. Controllability-aware unsuper-
520 vised skill discovery. *arXiv preprint arXiv:2302.05103*, 2023.
- 521 [49] Shubham Pateria, Budhitama Subagdja, Ah-Hwee Tan, and Chai Quek. End-to-end hierarchical
522 reinforcement learning with integrated subgoal discovery. *IEEE Transactions on Neural*
523 *Networks and Learning Systems*, 33(12):7778–7790, 2021.
- 524 [50] Judea Pearl. *Causality*. Cambridge university press, 2009.
- 525 [51] Silviu Pitis, Elliot Creager, and Animesh Garg. Counterfactual data augmentation using locally
526 factored dynamics. *Advances in Neural Information Processing Systems*, 33:3976–3990, 2020.
- 527 [52] Silviu Pitis, Elliot Creager, Ajay Mandlekar, and Animesh Garg. Mocoda: Model-based
528 counterfactual data augmentation. *Advances in Neural Information Processing Systems*, 35:
529 18143–18156, 2022.
- 530 [53] David Poole and Nevin Lianwen Zhang. Exploiting contextual independence in probabilistic
531 inference. *Journal of Artificial Intelligence Research*, 18:263–313, 2003.
- 532 [54] Rafael Rodriguez-Sanchez and George Konidaris. Learning abstract world models for value-
533 preserving planning with options. In *NeurIPS 2023 Workshop on Generalization in Planning*,
534 2023.
- 535 [55] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
536 policy optimization algorithms, 2017.
- 537 [56] Maximilian Seitzer, Bernhard Schölkopf, and Georg Martius. Causal influence detection for
538 improving efficiency in reinforcement learning. *Advances in Neural Information Processing*
539 *Systems*, 34:22905–22918, 2021.
- 540 [57] Wonil Song, Sangryul Jeon, Hyesong Choi, Kwanghoon Sohn, and Dongbo Min. Learning
541 disentangled skills for hierarchical reinforcement learning through trajectory autoencoder with
542 weak labels. *Expert Systems with Applications*, page 120625, 2023.
- 543 [58] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A
544 framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):
545 181–211, 1999.
- 546 [59] Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg,
547 David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning.
548 In *International Conference on Machine Learning*, pages 3540–3549. PMLR, 2017.

- 549 [60] Zizhao Wang, Jiaheng Hu, Peter Stone, and Roberto Martín-Martín. Elden: Exploration via
550 local dependencies. *Advances in Neural Information Processing Systems*, 36, 2024.
- 551 [61] Jiayi Weng, Huayu Chen, Dong Yan, Kaichao You, Alexis Duburcq, Minghao Zhang, Yi Su,
552 Hang Su, and Jun Zhu. Tianshou: A highly modularized deep reinforcement learning library.
553 *Journal of Machine Learning Research*, 23(267):1–6, 2022. URL [http://jmlr.org/papers/
554 v23/21-1127.html](http://jmlr.org/papers/v23/21-1127.html).
- 555 [62] Peter R Wurman, Samuel Barrett, Kenta Kawamoto, James MacGlashan, Kaushik Subramanian,
556 Thomas J Walsh, Roberto Capobianco, Alisa Devlic, Franziska Eckert, Florian Fuchs, et al.
557 Outracing champion gran turismo drivers with deep reinforcement learning. *Nature*, 602(7896):
558 223–228, 2022.
- 559 [63] Andrii Zadaianchuk, Maximilian Seitzer, and Georg Martius. Self-supervised visual reinforce-
560 ment learning with object-centric representations. In *International Conference on Learning
561 Representations*, 2021. URL <https://openreview.net/forum?id=xppLmXCb0w1>.
- 562 [64] Jesse Zhang, Haonan Yu, and Wei Xu. Hierarchical reinforcement learning by discovering
563 intrinsic options. *arXiv preprint arXiv:2101.06521*, 2021.

564 **A Factored Skills**

565 Learning to reach both a desired graph g and a diversity parameter b through primitive actions is
 566 challenging. First, different graphs often have substantially different characteristics, with some
 567 graphs that are easy to achieve (eg. action→agent), and others that are quite challenging and rare (eg.
 568 agent, knife, fruit→fruit). Not only would it be challenging for a single policy to encode all of these
 569 behaviors, the diversity parameter notwithstanding, but over-training the frequency at which certain
 570 graphs are called might vary significantly. Rather than trying to learn a single monolithic policy, then,
 571 we instead structure the skill parameterized policy π_{skill} as a collection of factored skills: $\pi_{\text{skill},i}$, for
 572 each factor $i \in \{1, \dots, N\}$.

573 This modification to the policy structure results in three changes: **1)** The upper-level action space
 574 passes a single row of the graph \mathcal{G} , denoted with g_i , and the desired factor i . **2)** Instead of achieving
 575 an entire graph use the achieved row $\mathbb{1}[g_{\text{achieved},i} = g_i]$. **3)** The history of seen graphs \mathcal{H} is replaced
 576 with a history of factored graph rows \mathcal{H}_f .

577 Define the history of graph rows as $\mathcal{H}_f := \{\text{unique}(i, g_{\text{achieved},i} \ \forall i \in 1, \dots, N \ \forall g_{\text{achieved}} \in \mathcal{D})\}$.
 578 This takes the unique graph rows from all those seen in previous data. Then the upper policy uses the
 579 same historical sampling procedure as with unfactored graphs: the policy samples discretely from
 580 the new history, which will by default return i, g_i , a graph row, and the desired factor. This resolves
 581 points **1,3**. Point **2** is addressed by replacing Equation 4 with $\mathbb{1}[g_{\text{achieved},i} = g_i]$.

582 Empirically, we found that without this change, the lower policy rarely learns anything, even simple
 583 control of the agent.

584 **B Environment Details**

585 In this section, we provide a detailed description of the environment, including its semantic stages
 586 representing internal progress toward task completion, state space, and action space. We also highlight
 587 that while each task consists of multiple semantic stages, agents do not have access to this information.

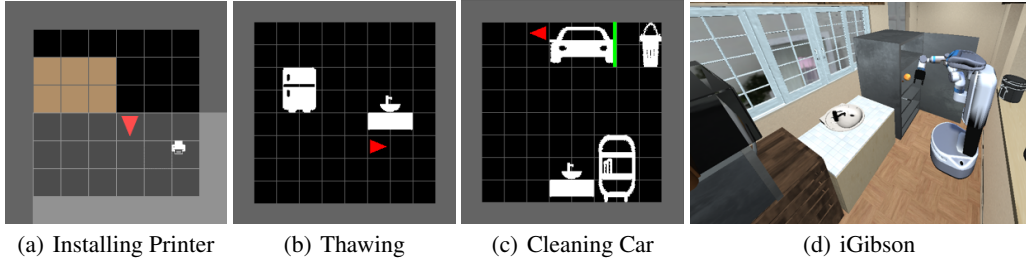


Figure 7: Environments.

588 **Installing Printer** As shown in Fig. 7(a), the Installing Printer environment is relatively simple,
 589 consisting of 3 factors: the agent, a printer, and a table. The task requires the agent to complete the
 590 following **stages**: (1) pick up the printer, (2) bring the printer to and place it on the table, and (3)
 591 turn on the printer. The discrete state space consists of (i) the agent’s position and direction, (ii) the
 592 positions of the printer and whether it is on or off, and (iii) the position of the table. The discrete
 593 action space consists of (i) moving forward, turning left or right, (ii) picking up / placing down the
 594 printer, and (iii) turning on / off the printer.

595 **Thawing** As shown in Fig. 7(b) and Fig. 8(a), the Thawing environment consists of 6 factors: the
 596 agent, a sink, a refrigerator, and three frozen objects: fish, olive, and date. Thawing each object
 597 requires the agent to complete the following **stages**: (1) move to and open the refrigerator, (2) take
 598 the frozen fish out of the refrigerator, (3) put the fish into the sink, and (4) turn on the sink to thaw
 599 it. The discrete state space consists of (i) the agent’s position and direction, (ii) the positions of all
 600 environment entities, (iii) whether the sink door is turned on, (iv) whether the refrigerator door is
 601 opened, and (v) the thawing status of three objects. The discrete action space consists of (i) moving
 602 forward, turning left or right, (ii) opening / closing the refrigerator, (iii) turning on / off the sink, and
 603 (iv) picking up / placing down each object.

604 **Cleaning Car** As shown in Fig. 7(c), the Cleaning Car environment consists of 7 factors: the agent,
 605 a car, a sink, a bucket, a shelf, a rag, and a piece of soap. Cleaning both the car and the rag requires
 606 the agent to complete the following **stages**: (1) take the rag off the shelf, (2) put it in the sink, (3)
 607 toggle the sink to soak the rag up, (4) clean the car with the soaked rag, (5) take the soap off the
 608 self, and (6) clean the rag with the soap inside the bucket. The discrete state space consists of (i) the
 609 agent’s position and direction, (ii) the positions of all environment entities, (iii) whether the sink is
 610 turned on, (iv) the soak status of the rag, (v) the cleanness of the rag, and (vi) the cleanness of the car.
 611 The discrete action space consists of (i) moving forward, turning left or right, (ii) turning on / off the
 612 sink, and (iii) picking up / placing down the rag / soap.

613 **iGibson** As shown in Fig. 7(d), the iGibson environment consists of 4 factors: the robot, a knife, a
 614 peach, and a sink. The robot can do the following things: (1) grasp peach: move close to the peach
 615 and grasp it, (3) wash peach: grasp the peach and place it into the sink, (3) grasp knife: move close to
 616 the knife and grasp it, (4) cut peach: grasp the knife and use it to cut the peach. The continuous state
 617 space consists of (i) the robot’s proprioception, (ii) the poses of all environment entities, and (iii)
 618 whether the peach is cut. The continuous action space consists of (i) end-effector position change, (ii)
 619 base linear and angular velocity, and (iii) gripper torque (to open/close the gripper).

620 C Local Dependencies and Causal Inference

621 In this work, we define local dependencies according to the state factors $X = (X^1, \dots, X^N)$
 622 and event of interest Y , which in the context of an MDP is a subset of the next state factors
 623 $X' = (X'^1, \dots, X'^N)$. In the factored MDP formulation [7], we assume that p , the transition
 624 dynamics, are represented by a dynamic Bayesian network (DBN) which is a time-directed bipartite
 625 graph, with edges only from factors in X to factors in X' . In this work, we assume that the underlying
 626 ground truth DBN, that is the transition function p , can be decomposed according to subsets of state
 627 factors \bar{X} , such there exists a $p^{\bar{X}}(Y = y | \bar{X} = x)$ for every state.

628 The factored transition dynamics analogizes with causal inference in the following way: If the state
 629 factors and next state factors are each assigned a causal variable by adding the assumption that they
 630 can be independently intervened on, and each next state variable carries an associated unobserved
 631 noise variable U^i , which we assume is independent of any X^k not connected to X'^j and any other
 632 next state variable X'^j , then we can represent the transition dynamics p with a structural causal model
 633 (SCM) [50], a graph connecting the causal variables in X to the causal variables in X' .

634 For a particular outcome variable Y that is one of the next state causal variables X' , we can describe
 635 local dependence in the RL context according to assumptions about the structural causal model.
 636 Represent the non-noise parents of Y as $\text{pa}(Y)$, and the noise parents as $\text{pa}_U(Y)$. Under normal
 637 causal assumptions, the structural causal model for Y is a function $f_Y(\text{pa}(Y), \text{pa}_U(Y)) = Y$. Define
 638 \bar{X} as a subset of the endogenous parents of Y and \bar{U} as an equivalent subset of the noise variables.
 639 Further define the values that $\text{pa}(Y)$, $\text{pa}_U(Y)$, \bar{X} , \bar{U} can take on as $\text{pa}(y)$, $\text{pa}_U(y)$, \bar{x} , \bar{u} respectively,
 640 and $(\text{pa}(\mathcal{Y}))$, $\bar{\mathcal{X}}$, $\bar{\mathcal{U}}$ as the set of states the parents of Y , the variables in \bar{X} and variables in \bar{U} can take
 641 on respectively.

642 To formalize local invariance, we add the assumption that f_Y can be decomposed into a series of func-
 643 tions $(f_{Y1}(\bar{X}_1 = \bar{x}_1, \bar{U}_1 = \bar{u}_1), \dots, f_{Yk}(\bar{X}_k = \bar{x}_k, \bar{U}_k = \bar{u}_k))$ and $g_Y(\text{pa}(Y) = \text{pa}(y), \text{pa}_U(Y) =$
 644 $\text{pa}_U(y))$, where each $f_{Yi} : \bar{\mathcal{X}} \times \bar{\mathcal{U}} \rightarrow \mathcal{Y}$ and $g : \text{pa}(\mathcal{Y}) \rightarrow \{1, \dots, k\}$, a function mapping the parents
 645 of Y to one of the functions. Then if f is represented as:

$$f(\text{pa}(x), \text{pa}_U(y)) := \sum_{i=1}^k \mathbb{1}(g_Y(\text{pa}(y), \text{pa}_U(y)) = i) f_{Yi}(\bar{x}_i, \bar{u}_i) \quad (6)$$

646 The local dependence of $Y = y$ in a particular state (x, x') is then the set of variables in \bar{X}_i for the
 647 particular i where $\mathbb{1}(g_Y(\text{pa}(y), \text{pa}_U(y)) = i) = 1$, and the pCMI test is a way of uncovering these
 648 local dependencies from observational data.

649 Local dependence has been investigated in the field of context-specific independence [53, 8], which
 650 seeks to find particular assignments of a subset of the causal variables under which an outcome
 651 is independent of some subset of the inputs. In particular, context-set specific independence [8]
 652 determines if a variable is independent of other variables on a particular subset of states, described as

Table 1: Parameters of Skill Learning and Task Learning. Parameters shared if not specified.

	Name	Environments			
		Printer	Thawing	Cleaning Car	iGibson
Skill Policy	algorithm		Rainbow		TD3
	n step		3		5
	skill horizon		30		100
	exploration noise		0.4		0.2
	optimizer			Adam	
	learning rate			3×10^{-4}	
	batch size		64		
Graph Selection Policy	algorithm			PPO	
	optimizer			Adam	
	learning rate			1×10^{-4}	
	batch size			1024	
	clip ratio			0.1	
	MLP size			[512, 512]	
	GAE λ			0.95	
entropy coefficient			0.1		
Learned Dynamics Model	optimizer			Adam	
	learning rate			3×10^{-4}	
	batch size			128	
	number of attention layers			1	
	attention embedding size			128	
	number of heads			4	
Task Skill Selection Policy	algorithm			PPO	
	optimizer			Adam	
	learning rate			1×10^{-4}	
	batch size			1024	
	clip ratio			0.1	
	MLP size			[512, 512]	
	GAE λ			0.95	
	entropy coefficient			0.02	
Training	# of random seeds			5	
	diversity reward coefficient β			0.5	

653 the partial context set. While our work uses the pCMI test described in Equation 3, context-specific
654 independence focuses on complete independence using knowledge of the structural model.

655 Alternatively, interactions can be viewed as the causes (\bar{X}) of particular effects (Y), which have
656 also been investigated under the description of token or actual cause [24] (as opposed to general
657 cause). Actual cause utilizes a series of counterfactual tests to determine if a cause is necessary,
658 sufficient, and minimal for an outcome. Actual cause has primarily been applied in simple, discrete
659 examples [4, 23], making it difficult to directly apply to RL. However, recent work has incorporated
660 the notion of context-specific independence and extended actual cause to more complex domains [14].

661 D Implementation Details

662 The hyperparameters of skill learning and task learning can be found in Table 1. As it is challenging
663 to identify local dependencies using learned dynamics models in Thawing and iGibson environments,
664 we use ground truth local dependencies from simulator. The codebase is built on tianshou [61] for
665 backend RL, though with significant modifications.

666 The 5 seeds selected are 0 - 4. The experiments were conducted on machines of the following
667 configurations:

- 668 • Nvidia A40 GPU; Intel(R) Xeon(R) Gold 6342 CPU @2.80GHz
- 669 • Nvidia A100 GPU; Intel(R) Xeon(R) Gold 6342 CPU @2.80GHz

670 **E Skill Visualizations**

671 In Figure 8 we visualize three challenging long-horizon skills learned by SkiLD: thawing the olive,
672 cleaning the car, and cutting the peach. All of these skills require a sequence of interactions that
673 is difficult to recover without directed behavior. Thus, comparable baselines do not learn skills of
674 similar complexity. More skill visualizations can be found at: <https://sites.google.com/view/skild>.

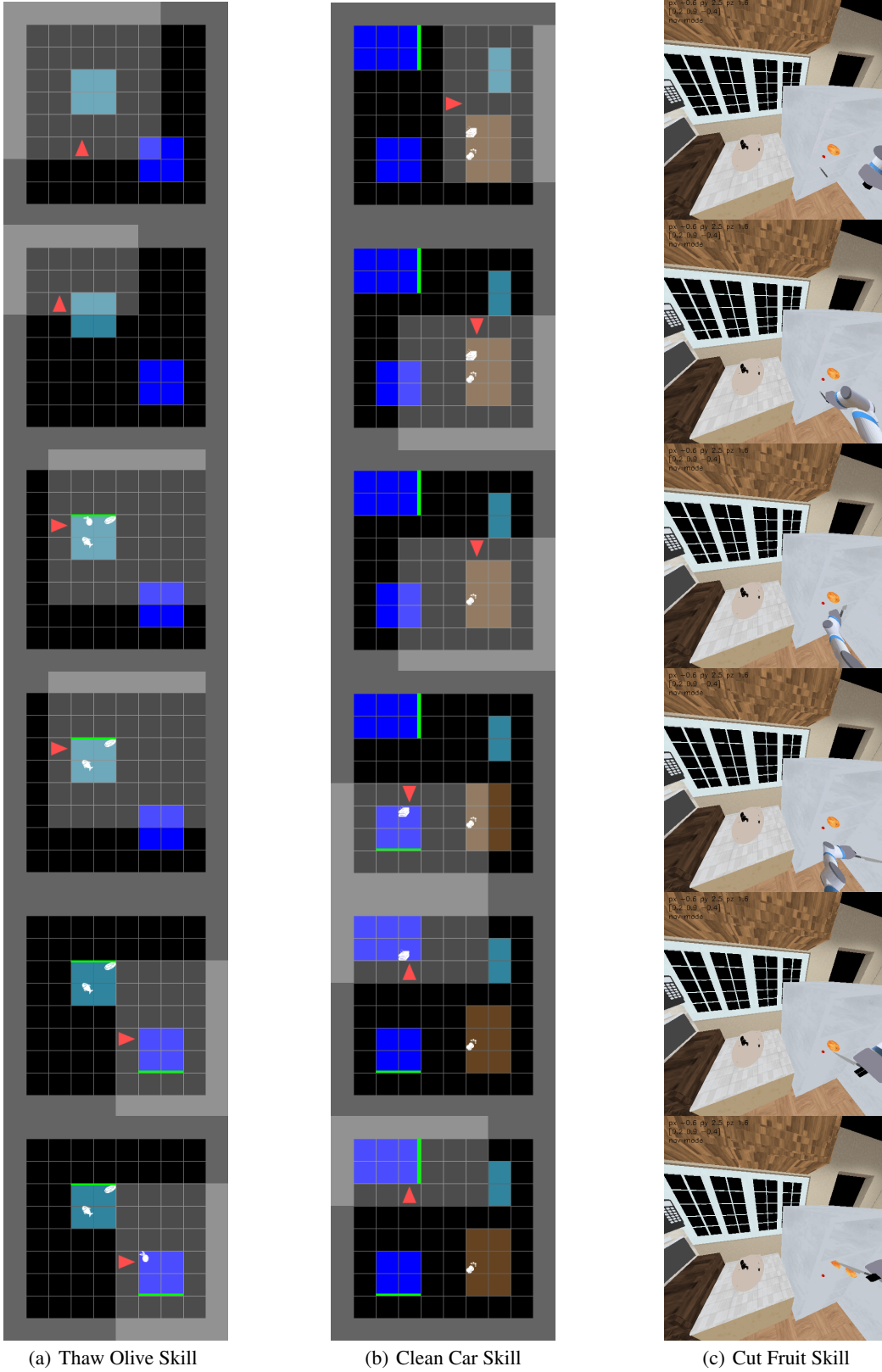


Figure 8: Policy rollouts for learned policies that achieve long horizon tasks **(a)** Mini-BH thaw olive, **(b)** Mini-BH clean car, **(b)** iGibson cut peach.