

---

# Consciousness-Inspired Spatio-Temporal Abstractions for Better Generalization in Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Inspired by human conscious planning, we propose **Skipper**, a model-based rein-  
2 forcement learning framework utilizing spatio-temporal abstractions to generalize  
3 better in novel situations. It automatically decomposes the given task into smaller,  
4 more manageable subtasks, and thus enables sparse decision-making and focused  
5 computation on the relevant parts of the environment. The decomposition relies  
6 on the extraction of an abstracted proxy problem represented as a directed graph,  
7 in which vertices and edges are learned end-to-end from hindsight. Our theoret-  
8 ical analyses provide performance guarantees under appropriate assumptions and  
9 establish where our approach is expected to be helpful. Generalization-focused  
10 experiments validate **Skipper**'s significant advantage in zero-shot generalization,  
11 compared to some existing state-of-the-art hierarchical planning methods.

## 12 1 Introduction

13 Attending to relevant aspects in both time and space, human conscious planning breaks down long-  
14 horizon tasks into more manageable steps, each of which can be narrowed down further. Stemming  
15 from consciousness in the first sense (C1) [15], this type of planning focuses attention on mostly the  
16 important decision points [63] and relevant environmental factors linking the decision points [66],  
17 thus *operating abstractly both in time and in space*. In contrast, existing Reinforcement Learning (RL)  
18 agents either operate solely based on intuition (model-free methods) or are limited to reasoning over  
19 mostly relatively shortsighted plans (model-based methods) [29]. The intrinsic limitations constrain  
20 the real-world application of RL under a glass ceiling formed by challenges of generalization.

21 Building on our previous work on conscious planning [73], we take inspirations to develop a planning  
22 agent that automatically decomposes the complex task at hand into smaller subtasks, by constructing  
23 abstract “proxy” problems. A proxy problem is represented as a graph where 1) the vertices consist  
24 of states imagined by a generative model, corresponding to sparse decision points; and 2) the edges,  
25 which define temporally-extended transitions, are constructed by focusing on a small amount of  
26 relevant information from the states, using an attention mechanism. Once a proxy problem is  
27 constructed and the agent solves it to form a plan, each of the edges defines a new sub-problem,  
28 on which the agent will focus next. This divide-and-conquer strategy allows constructing partial  
29 solutions that generalize better to new situations, while also giving the agent flexibility to construct  
30 abstractions necessary for the problem at hand. Our theoretical analyses establish guarantees on the  
31 quality of the solution to the overall problem.

32 We also examine empirically advantages of out-of-training-distribution generalization of our method  
33 after using only a few training tasks. We show through detailed controlled experiments that the  
34 proposed framework, named **Skipper**, in most cases performs significantly better in terms of zero-shot  
35 generalization, compared to the baselines and to some state-of-the-art Hierarchical Planning (HP)  
36 methods [45, 23].

## 37 2 Preliminaries

38 **Reinforcement Learning & Problem Setting.** An RL agent interacts with an environment via a  
39 sequence of actions to maximize its cumulative reward. The interaction is usually modeled as a  
40 Markov Decision Process (MDP)  $\mathcal{M} \equiv \langle \mathcal{S}, \mathcal{A}, P, R, d, \gamma \rangle$ , where  $\mathcal{S}$  and  $\mathcal{A}$  are the set of states and  
41 actions,  $P : \mathcal{S} \times \mathcal{A} \rightarrow \text{Dist}(\mathcal{S})$  is the state transition function,  $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is the reward  
42 function,  $d : \mathcal{S} \rightarrow \text{Dist}(\mathcal{S})$  is the initial state distribution, and  $\gamma \in [0, 1]$  is the discount factor. The  
43 agent needs to learn a policy  $\pi : \mathcal{S} \rightarrow \text{Dist}(\mathcal{A})$  that maximizes the value of the states, *i.e.* the expected  
44 discounted cumulative reward  $\mathbb{E}_{\pi, P}[\sum_{t=0}^{T_{\perp}} \gamma^t R(S_t, A_t, S_{t+1}) | S_0 \sim d]$ , where  $T_{\perp}$  denotes the time  
45 step at which the episode terminates. A value estimator  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  can be used to search for  
46 a good policy. However, real-world problems can be partially observable, meaning that, instead of  
47 states, the agent receives an observation  $x_{t+1} \in \mathcal{X}$ , where  $\mathcal{X}$  is the observation space. The agent  
48 needs to infer the state from the sequence of observations, usually with a state encoder.

49 One important goal of RL is to achieve high (generalization) performance on evaluation tasks after  
50 learning from a limited number of training tasks, where the evaluation and training distributions may  
51 differ; for instance, a policy for a robot may need to be trained in a simulated environment for safety  
52 reasons, but would need to be deployed on a physical device, a setting called *sim2real*. Discrepancy  
53 between task distributions is often recognized as a major reason why RL agents are yet to be applied  
54 pervasively in the real world [28]. To address this issue, in this paper, agents are trained on a small  
55 set of fixed training tasks, then evaluated in unseen tasks, where there are environmental variations,  
56 but the core strategies needed to finish the task remain consistent. To generalize well, the agents need  
57 to build learned skills which capture the consistent knowledge across tasks.

58 **Deep Model-based RL.** Deep model-based RL uses predictive or generative models to help search  
59 for policies [59]. For generalization, rich models, expressed by Neural Networks (NNs), may capture  
60 generalizable information and infer latent causal structure. *Background* planning agents *e.g.*, Dreamer  
61 [25] use a model as a data generator to improve the value estimators and policies, which executes  
62 in background without directly engaging in the environment [61]. These agents do not improve on  
63 the trained policy at decision time. In contrast, *decision-time* planning agents *e.g.*, MuZero [54]  
64 and PlaNet [24] actively use models to make better decisions. Recently, [1] suggests that the latter  
65 approach provides better generalization, aligning with observations from cognitive behaviors [40].

66 **Options & Goal-Conditioned RL.** Temporal abstraction allows agents to use sub-policies, and  
67 to model the environment over extended time scales, to achieve both better generalization and the  
68 divide and conquer of larger problems. Options and their models provide a formalism for temporal  
69 abstraction in RL [63]. Each option consists of an initiation condition, a policy, and a termination  
70 condition. For any set of options defined on an MDP, the decision process that selects only among  
71 those options, executing each to termination, is a Semi-MDP (SMDP) [63, 49], consisting of the  
72 set of states  $\mathcal{S}$ , the set of options  $\mathcal{O}$ , and for each state-option pair, an expected return, and a joint  
73 distribution of the next state and transit time. In this paper, we focus on goal-conditioned options,  
74 where the initiation set covers the whole state space  $\mathcal{S}$ . Each such option is a tuple  $o = \langle \pi, \beta \rangle$ , where  
75  $\pi : \mathcal{S} \rightarrow \text{Dist}(\mathcal{A})$  is the (intra-)option policy and  $\beta : \mathcal{S} \rightarrow \{0, 1\}$  indicates when a goal state is  
76 reached. Hindsight Experience Replay (HER) [3] is often used to train goal-conditioned options by  
77 sampling a transition  $\langle x_t, a_t, r_{t+1}, x_{t+1} \rangle$  together with an additional observation  $x^{\circ}$  from the same  
78 trajectory, which is re-labelled as a “goal”.

## 79 3 Skipper: Spatially & Temporally Abstract Planning

80 In this section, we describe the main ingredients of **Skipper** - a framework that formulates a **proxy**  
81 problem for a given task, solves this problem, and then proceeds to “fill in” the details of the plan.

### 82 3.1 Proxy Problems

83 Proxy problems are finite graphs constructed at decision-time, whose vertices are states and whose  
84 directed edges estimate transitions between the vertices, as shown in Fig. 1. We call the states  
85 selected to be vertices of the proxy problems *checkpoints*, to differentiate from other uninvolved  
86 states. The current state is always included as one of the vertices. The checkpoints are proposed by a  
87 generative model and represent some states that the agent might experience in the current episode,  
88 often denoted as  $S^{\circ}$  in this paper. Each edge is annotated with estimates of the cumulative discount

89 and reward associated with the transition between the connected checkpoints; these estimates are  
 90 learned over the **relevant** aspects of the environment and **depend** on the agent’s capability. As the  
 91 low-level policy implementing checkpoint transitions improves, the edges strengthen. Planning in a  
 92 proxy problem is temporally abstract, since the checkpoints act as sparse decision points. Estimating  
 93 each checkpoint transition is spatially abstract, as an option corresponding to such a task would base  
 94 its decisions only on some aspects of the environment state [7, 34], to improve generalization as well  
 95 as computational efficiency [73].

96 A proxy problem can be viewed as a deterministic SMDP, where each  
 97 directed edge is implemented as a checkpoint-conditioned option. It can  
 98 be fully described by the discount and reward matrices,  $\Gamma^\pi$  and  $V^\pi$ , where  
 99  $\gamma_{ij}^\pi$  and  $v_{ij}^\pi$  are defined as:

$$\gamma_{ij}^\pi \doteq \mathbb{E}_\pi [\gamma^{T_\perp} | S_0 = s_i, S_{T_\perp} = s_j] \quad (1)$$

$$v_{ij}^\pi \doteq \mathbb{E}_\pi \left[ \sum_{t=0}^{T_\perp} \gamma^t R_t | S_0 = s_i, S_{T_\perp} = s_j \right]. \quad (2)$$

100 By planning with  $\Gamma^\pi$  and  $V^\pi$ , e.g. using SMDP value iteration [63], we  
 101 can solve the proxy problem, and form a jumpy plan to travel between  
 102 states in the original problem. If the proxy problems can be estimated  
 103 well, the obtained solution will be of good quality, as established in the  
 104 following theorem:

105 **Theorem 1** Let  $\mu$  be the SMDP policy (high-level) and  $\pi$  be the low-level  
 106 policy. Let  $\hat{V}^\pi$  and  $\hat{\Gamma}^\pi$  denote learned estimates of the SMDP model. If  
 107 the estimation accuracy satisfies:

$$\begin{aligned} |v_{ij}^\pi - \hat{v}_{ij}^\pi| < \epsilon_v v_{max} \ll (1 - \gamma) v_{max} & \quad \text{and} & (3) \\ |\gamma_{ij}^\pi - \hat{\gamma}_{ij}^\pi| < \epsilon_\gamma \ll (1 - \gamma)^2 & \quad \forall i, j. \end{aligned}$$

108 Then, the estimated value of the composite  $\hat{v}_{\mu \circ \pi}(s)$  is accurate up to error terms linear in  $\epsilon_v$  and  $\epsilon_\gamma$ :

$$\hat{v}_{\mu \circ \pi}(s) \doteq \sum_{k=0}^{\infty} \hat{v}_\pi(s_k^\circ | s_{k+1}^\circ) \prod_{\ell=0}^{k-1} \hat{\gamma}_\pi(s_\ell^\circ | s_{\ell+1}^\circ) = v_{\mu \circ \pi}(s) \pm \frac{\epsilon_v v_{max}}{1 - \gamma} \pm \frac{\epsilon_\gamma v_{max}}{(1 - \gamma)^2} + o(\epsilon_v + \epsilon_\gamma)$$

109 where  $\hat{v}_\pi(s_i | s_j) \equiv \hat{v}_{ij}^\pi$  and  $\hat{\gamma}_\pi(s_i | s_j) \equiv \hat{\gamma}_{ij}^\pi$ , and  $v_{max}$  denotes the maximum value.

110 The theorem indicates that once the agent achieves high accuracy estimation of the model for the  
 111 proxy problem and a near-optimal lower-level policy  $\pi$ , it converges toward optimal performance  
 112 (proof in Appendix D.2). The theorem also makes no assumption on  $\pi$ , since it would likely be  
 113 difficult to learn a good  $\pi$  for far away targets. Despite the theorem’s generality, in the experiments,  
 114 we limit ourselves to navigation tasks with sparse rewards for reaching goals, where the goals are  
 115 included as permanent vertices in the proxy problems. This is a case where the accuracy assumption  
 116 can be met non-trivially, i.e., while avoiding degenerate proxy problems whose edges involve no  
 117 rewards. Following Thm. 1, we train estimators for  $v_\pi$  and  $\gamma_\pi$  and refer to this as *edge estimation*.

### 118 3.2 Design Choices

119 To implement planning over proxy problems, **Skippier** embraces the following design choices:

120 **Decision-time planning** is employed due to its ability to improve the policy in novel situations;

121 **Spatio-temporal abstraction**: temporal abstraction breaks down the given task into smaller ones,  
 122 while spatial abstraction<sup>1</sup> over the state features improves local learning and generalization;

123 **Higher quality proxies**: we introduce pruning techniques to improve the quality of proxy problems;

124 **Learning end-to-end from hindsight, off-policy**: to maximize sample efficiency and the ease of  
 125 training, we propose to use auxiliary (off-)policy methods for edge estimation, and learn a context-  
 126 conditioned checkpoint generation, both from hindsight experience replay;

<sup>1</sup>We use “spatial abstraction” to denote specifically the behavior of constraining decision-making to the relevant environmental factors during an option. Please check Section 4 for discussions and more details.

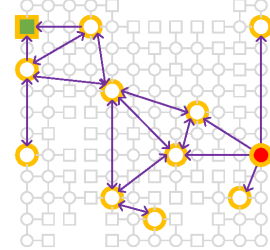


Figure 1: **A Proxy Problem on a Navigation Task**: the MDP of the original problem is in gray and the terminal states are marked with squares. An agent needs to get from the red position, to the goal (green). Distant goals can be reached by leveraging a proxy problem with 12 checkpoints (outlined orange).

127 **Delusion suppression:** we propose a delusion suppression technique to minimize the behavior of  
 128 chasing non-existent outcomes. This is done by exposing edge estimation to imagined targets that  
 129 would otherwise not exist in experience.

### 130 3.3 Problem 1: Edge Estimation

131 First, we discuss how to estimate the edges of the proxy problem, given a set of already generated  
 132 checkpoints. Inspired by conscious information processing in brains [15] and existing approach  
 133 in [64], we introduce a local perceptive field selector,  $\sigma$ , consisting of an attention bottleneck that  
 134 (soft-)selects the top- $k$  local segments of the full state (*e.g.* a feature map by a typical convolutional  
 135 encoder); all segments of the state compete for the  $k$  attention slots, *i.e.* irrelevant aspects of state  
 136 are discouraged or discarded, to form a partial state representation [41, 66, 73, 2]. We provide an  
 137 example in Fig. 2 (see purple parts). Through  $\sigma$ , the auxiliary estimators, to be discussed soon,  
 138 force the bottleneck mechanism to promote aspects relevant to the local estimation of connections  
 139 between the checkpoints. The rewards and discounts are then estimated from the partial state  $\sigma(S)$ ,  
 140 conditioned on the agent’s policy.

#### 141 3.3.1 Basis for Connections: Checkpoint-Achieving Policy

142 The low-level policy  $\pi$  maximizes an intrinsic reward, *s.t.* the target checkpoint  $S^\odot$  can be reached.  
 143 The choice of intrinsic reward is flexible; for example, one could use a reward of +1 when  $S_{t+1}$  is  
 144 within a small radius of  $S^\odot$  according to some distance metric, or use reward-respecting intrinsic  
 145 rewards that enable more sophisticated behaviors, as in [62]. In the following, for simplicity, we will  
 146 denote the checkpoint-achievement condition with equality:  $S_{t+1} = S^\odot$ .

#### 147 3.3.2 Estimate Connections

148 We learn the connection estimates with auxiliary reward signals that are designed to be not task-  
 149 specific [74]. These estimates are learned using C51-style distributional RL, where the output of each  
 150 estimator takes the form of a histogram over scalar support [14].

151 **Cumulative Reward.** The cumulative discounted task reward  $v_{ij}^\pi$  is learned by policy evaluation on  
 152 an auxiliary reward that is the same as the original task reward everywhere except when reaching  
 153 the target. Given a hindsight sample  $\langle x_t, a_t, r_{t+1}, x_{t+1}, x^\odot \rangle$  and the corresponding encoded sample  
 154  $\langle s_t, a_t, r_{t+1}, s_{t+1}, s^\odot \rangle$ , we train  $V_\pi$  with KL-divergence as follows:

$$\hat{v}_\pi(\sigma(s_t), a_t | \sigma(s^\odot)) \leftarrow \begin{cases} R(s_t, a_t, s_{t+1}) + \gamma \hat{v}_\pi(\sigma(s_{t+1}), a_{t+1} | \sigma(s^\odot)) & \text{if } s_{t+1} \neq s^\odot \\ R(s_t, a_t, s_{t+1}) & \text{if } s_{t+1} = s^\odot \end{cases} \quad (4)$$

155 where  $\sigma(s)$  is the spatially-abstracted from the full state  $s$  and  $a_{t+1} \sim \pi(\cdot | \sigma(s_{t+1}), \sigma(s^\odot))$ .

156 **Cumulative Distances / Discounts.** With C51 and uniform supports, the cumulative discount leading  
 157 to  $s^\odot$  under  $\pi$  is unfortunately more difficult to learn than  $V_\pi$ , since the prediction would be heavily  
 158 skewed towards 1 if  $\gamma \approx 1$ . Yet, we can instead effectively estimate cumulative (truncated) distances  
 159 (or trajectory length) under  $\pi$ . Such distances can be learned with policy evaluation, where the  
 160 auxiliary reward is +1 on every transition, except at the targets:

$$D_\pi(\sigma(s_t), a_t | \sigma(s^\odot)) \leftarrow \begin{cases} 1 + D_\pi(\sigma(s_{t+1}), a_{t+1} | \sigma(s^\odot)) & \text{if } s_{t+1} \neq s^\odot \\ 1 & \text{if } s_{t+1} = s^\odot \\ \infty & \text{if } s_{t+1} \text{ is terminal and } s_{t+1} \neq s^\odot \end{cases}$$

161 where  $a_{t+1} \sim \pi(\cdot | \sigma(s_{t+1}), \sigma(s^\odot))$ . The cumulative discount is then recovered by replacing the  
 162 support of the output distance histogram with the corresponding discounts. Additionally, the learned  
 163 distance is used to prune unwanted checkpoints to simplify the proxy problem, as well as prune  
 164 far-fetched edges. The details of pruning will be presented shortly.

166 Please refer to the Appendix D.1 for the properties of the learning rules for  $\hat{v}_\pi$  and  $\hat{\gamma}_\pi$ .

### 167 3.4 Problem 2: Vertex Generation

168 The checkpoint generator aims to directly imagine the possible future states *without needing to know*  
 169 *how exactly the agent might reach them nor worrying about if they are reachable*. The feasibility of  
 170 checkpoint transitions will be abstracted by the connection estimates instead.

171 To make the checkpoint generator generalize well across diverse tasks, while still being able to  
 172 capture the underlying causal mechanisms in the environment (a challenge for existing model-based  
 173 methods) [71], we propose that the checkpoint generator learns to split the state representation into  
 174 two parts: an episodic context and a partial description. In a navigation problem, for example, as in  
 175 Fig. 2, a context could be a representation of the map of a gridworld, and the partial description be  
 176 the 2D-coordinates of the agent’s location. In different contexts, the same partial description could  
 177 correspond to very different states. Yet, within the same context, we should be able to recover the  
 178 same state given the same partial description.

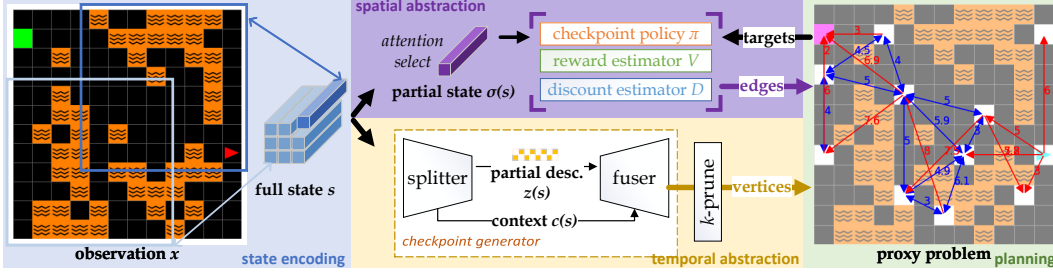


Figure 2: **Skipper Framework:** 1) Partial states consist of a few local fields, soft-selected via top- $k$  attention [22]. **Skipper**’s edge estimations and low-level behaviors  $\pi$  are based on the partial states. 2) The checkpoint generator learns by splitting the full state into context and partial descriptions, and fusing them to reconstruct the input. It imagines checkpoints by sampling partial descriptions and combining them with the episodic contexts; 3) We prune the vertices and edges of the denser graphs to extract sparser proxy problems. Once a plan is formed, the immediate checkpoint target is used to condition the policy. In the proxy problem example, blue edges are estimated to be bidirectional and red edges have the other direction pruned.

179 As shown in Fig. 2, this information split is achieved using two functions: the *splitter*  $\mathcal{E}_{CZ}$ , which  
 180 maps the input state  $S$  into a representation of a context  $c(S)$  and a partial description  $z(S)$ , as well  
 181 as the *fuser*  $\oplus$  which, when applied to the input  $\langle c, z \rangle$ , recovers  $S$ . In order to achieve consistent  
 182 context extraction across states in the same episode, at training time, we force the context to be  
 183 extracted from other states in the same episode, instead of the input.

184 We sample in hindsight a diverse distribution of target encoded (full) states  $S^\odot$ , given any current  
 185  $S_t$ . Hence, we make the generator a conditional Variational AutoEncoder (VAE) [60] which learns a  
 186 distribution  $p(S^\odot|C(S_t)) = \sum_z p(S^\odot|C(S_t), z)p(z|C(S_t))$ , where  $C(S_t)$  is the extracted context  
 187 from  $S_t$  and  $z$ s are the partial descriptions. We train the generator by minimizing the evidence lower  
 188 bound on  $\langle S_t, S^\odot \rangle$  pairs chosen with HER.

189 Similarly to [25], we constrain the partial description as a bundle of binary variables and train them  
 190 with the straight-through gradients [8]. These binary latents can be easily sampled or composed for  
 191 generation. Compared to models such as that in Director [23], which generates intermediate goals  
 192 given the on-policy trajectory, ours can generate and handle a more diverse distribution of states,  
 193 beneficial for planning in novel scenarios.

### 194 3.4.1 Pruning

195 In this paper, we limit ourselves only to checkpoints from a return-unaware conditional generation  
 196 model, leaving the question of how to improve the quality of the generated checkpoints for future  
 197 work. Without learning, the proxy problem can be improved by making it more sparse, and making  
 198 the proxy problem vertices more evenly spread in state space. To achieve this, we propose a pruning  
 199 algorithm based on  $k$ -medoids clustering [30], which only requires pairwise distance estimates  
 200 between states. During proxy problem construction, we first sample a larger number of checkpoints,  
 201 and then cluster them and select the centers (which are always real states instead of imaginary  
 202 weighted sums of state representations).

203 Notably, for sparse reward tasks, the generator cannot guarantee the presence of the rewarding  
 204 checkpoints in the proposed proxy problem. We could remedy this by explicitly learning the  
 205 generation of the rewarding states with another conditional generator. These rewarding states should  
 206 be kept as vertices (immune from pruning).

207 In addition to pruning the vertices, we also prune the edges according to a distance threshold, *i.e.*, all  
 208 edges with estimated distance over the threshold are deleted from the complete graph of the pruned

209 vertices. This biases potential plans towards shorter-length, smaller-scale sub-problems, as far-away  
210 checkpoints are difficult for  $\pi$  to achieve, trading optimality for robustness.

### 211 3.4.2 Safety & Delusion Control

212 Model-based HRL agents can be prone to blindly optimizing for objectives without understanding the  
213 consequences [36, 46]. We propose a technique to suppress delusions by exposing edge estimation to  
214 potentially delusional targets that do not exist in the experience replay buffer. Details and examples  
215 are provided in the Appendix.

## 216 4 Related Works & Discussions

217 **Temporal Abstraction.** Resembling attention, choosing a checkpoint target is a selection towards  
218 certain decision points in the dimension of time, *i.e.* a form of temporal abstraction. Constraining  
219 options, **Skipper** learns the options targeting certain “outcomes”, which dodges the difficulties  
220 of option collapse [5] and option outcome modelling by design. The constraints indeed shift the  
221 difficulties to generator learning [58, 65]. We expect this to entail benefits where states are easy to  
222 learn and generate, and / or in stochastic environments where the outcomes of unconstrained options  
223 are difficult to learn. Constraining options was also investigated in [56] in an unsupervised setting.

224 **Spatial Abstraction** is different from “state abstraction” [52, 33], which evolved to be a synonym for  
225 state space partitioning [37]. Spatial abstraction, defined to capture the behavior of conscious planning  
226 in the spatial dimension, focuses on the **within-state** partial selection of the environmental state for  
227 decision-making. It corresponds naturally to the intuition that state representations should contain  
228 useful aspects of the environment, while not all aspects are useful for a particular intent. Efforts  
229 toward spatial abstraction are traceable to early hand-coded proof-of-concepts proposed in *e.g.* [16].  
230 Until only recently, attention mechanisms had primarily been used to construct state representations  
231 in model-free agents for sample efficiency purposes, without the focus on generalization [41, 38, 66].  
232 In [20, 70, 55], 3 more recent model-based approaches, spatial abstractions are attempted to remove  
233 visual distractors. Concurrently, emphasizing on generalization, our previous work [73] used spatially-  
234 abstract partial states in decision-time planning. We proposed an attention bottleneck to dynamically  
235 select a subset of environmental entities during the atomic-step forward simulation, without explicit  
236 goals provided as in [70]. **Skipper**’s checkpoint transition is a step-up from our old approach, where  
237 we now show that spatial abstraction, an overlooked missing flavor, is as crucial for longer-term  
238 planning as temporal abstraction [34].

239 **Task Abstraction via Goal Composition** The early work [39] suggested to use bottleneck states  
240 as subgoals to abstract given tasks into manageable steps. [43, 19] use generative model to imagine  
241 subgoals while [18] search directly on the experience replay. In [31], promising states to explore  
242 are generated and selected with shortest-path algorithms. Similar ideas have been attempted for  
243 guided exploration [17, 35]. Similar to [23], [13] generate fixed-steps ahead subgoals for reasoning  
244 tasks, while [6] augments the search graph by states reached fixed-steps ahead. [45, 69, 57] employ  
245 CEM to plan a chain of subgoals towards the task goal [50]. **Skipper** utilizes proxy problems to  
246 abstract the given tasks via spatio-temporal abstractions [6]. Checkpoints can be seen as sub-goals  
247 that generalize the notion of “landmarks” or “waypoints” in [63, 16, 53]. [72] used latent landmark  
248 graphs as high-level guidance, where the landmarks are sparsified with weighted sums in the latent  
249 space to compose subgoals. In comparison, our checkpoint pruning selects a subset of generated  
250 states, which is less prone to issues created by weighted sums.

251 **Planning Estimates.** [72] propose a distance estimate with an explicit regression. With TDMs  
252 [48], LEAP [45] embraces a sparse intrinsic reward based on distances to the goal. Contrasting with  
253 our distance estimates, there is no empirical evidence of TDMs’ compatibility with stochasticity  
254 and terminal states. Notably, [18] employs a similar distance learning scheme to learn the shortest  
255 path distance between states found in the experience replay; while our estimators learn the distance  
256 conditioned on evolving policies. Such aspect was also investigated in [42].

257 **Decision-Time HP** with evolutionary algorithms were investigated in [44, 24, 45].

## 258 5 Experiments

259 As introduced in Sec. 2, our first goal is to test the zero-shot generalization ability of trained agents.  
260 To fully understand the results, it is necessary to have precise control of the difficulty of the training  
261 and evaluation tasks. Also, to validate if the empirical performance of our agents matches the formal  
262 analyses (Thm. 1), we need to know how close to the (optimal) ground truth our edge estimations  
263 and checkpoint policies are. These goals lead to the need for environments whose ground truth  
264 information (optimal policies, true distances between checkpoints, etc) can be computed. Thus,  
265 we base our experimental setting on the MiniGrid-BabyAI framework [10, 9, 27]. Specifically, we  
266 build on the experiments used in our previous works [73, 1]: the agent needs to navigate to the goal  
267 from its initial state in gridworlds filled with terminal lava traps generated randomly according to a  
268 difficulty parameter, which controls their density. During evaluation, the agent is always spawned  
269 at the opposite side from the goals. During training, the agent’s position is uniformly initialized to  
270 speed up training. We provide results for non-uniform training initialization in the Appendix.

271 These fully observable tasks prioritize on the challenge of reasoning over causal mechanisms over  
272 learning representations from complicated observations. Across all experiments, we sample training  
273 tasks from an environment distribution of difficulty 0.4: each cell in the field has probability 0.4 to  
274 be filled with lava while guaranteeing a path from the initial position to the goal. The evaluation  
275 tasks are sampled from a gradient of OOD difficulties - 0.25, 0.35, 0.45 and 0.55, where the training  
276 difficulty acts as mean. To step up the long(er) term generalization difficulty compared to existing  
277 work, we conduct experiments done on large,  $12 \times 12$  maze sizes, (see the visualization in Fig 2).  
278 The agents are trained for  $1.5 \times 10^6$  interactions. The compared agents include:

279 **Skipper-once**: A **Skipper** agent that generates one proxy problem at the start of the episode, and the  
280 replanning (choosing a checkpoint target based on the existing proxy problem) only triggers a quick  
281 re-selection of the immediate checkpoint target;

282 **Skipper-regen**: A **Skipper** agent that re-generates a proxy problem when replanning is triggered;

283 **modelfree**: A model-free baseline agent sharing the same base architecture with the **Skipper** variants  
284 - a prioritized distributional Double DQN [14, 68];

285 **Director**: A tuned Director agent [23] fed with simplified visual inputs. Since Director discards  
286 trajectories that are not long enough for training purposes, we make sure that the same amount of  
287 training data is gathered as for the other agents;

288 **LEAP**: A re-implemented LEAP for discrete action spaces. Due to low performance, we replaced the  
289 VAE and the distance learning mechanisms with our counterparts. We waived the interaction costs  
290 for its generator pretraining stage, only showing the second stage of RL pretraining.

291 Please refer to the Appendix for more details and insights on these agents.

### 292 5.1 Generalization Performance

293 Fig. 3 shows how the agents’ generalization performance evolves during training. These results  
294 are obtained with 50 fixed sampled training tasks (different 50s for each seed), a representative  
295 configuration of different numbers of training tasks including  $\{1, 5, 25, 50, 100, \infty\}^2$ , whose results  
296 are in the Appendix. In Fig. 3 a), we observe how well an agent performs on its training tasks. If an  
297 agent performs well here but badly in b), c), d) and e), *e.g.* the **modelfree** baseline, then we suspect  
298 that it overfitted on training tasks, likely indicating a reliance on memorization [11].

299 We observe a (statistically-)significant advantage in the generalization performance of the **Skipper**  
300 agents throughout training. We have also included significance tests and power analyses [12, 47]  
301 in the Appendix, together with results for other training configurations. The **regen** variant exhibits  
302 dominating performance over all others. This is likely due to the frequent reconstruction of the  
303 graph makes the agent less prone to being trapped in a low-quality proxy problem and provides extra  
304 adaptability in novel scenarios (more discussions in the Appendix). During training, **Skipper**s behave  
305 less optimally than expected, despite the strong generalization on evaluation tasks. As our ablation  
306 results and theoretical analyses consistently show, such a phenomenon is a composite outcome  
307 of inaccuracies both in the proxy problem and the checkpoint policy. One major symptom of an

---

<sup>2</sup> $\infty$  training tasks mean that an agent is trained on a different task for each episode. In reality, this may lead to prohibitive costs in creating the training environment.

308 inaccurate proxy problem is that the agent would chase delusional targets. We address this behavior  
 309 with the delusion suppression technique, to be discussed in the Appendix.

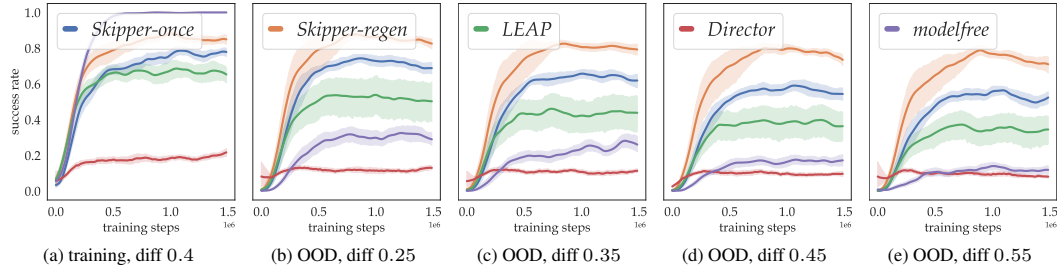


Figure 3: **Generalization Performance of Agents During Training:** the  $x$ -axes correspond to training progress, while the aligned  $y$ -axes represent the success rate of episodes (optimal is 1.0). Each agent is trained with 50 tasks. Each data point is the average success rate over 20 evaluation episodes, and each error bar (95% confidence interval) is processed from 20 independent seed runs. Training tasks performance is shown in (a) while OOD evaluation performance is shown in (b), (c), (d), (e).

310 Better than the **modelfree** baseline, LEAP obtains reasonable generalization performance, despite the  
 311 extra budget it needs for pretraining. In the Appendix, we show that LEAP benefits largely from the  
 312 delusion suppression technique. This indicates that optimizing for a path in the latent space may be  
 313 prone to errors caused by delusional subgoals. Lastly, we see that the Director agents suffer in these  
 314 experiments despite their good performance in the single environment experimental settings reported  
 315 by [23]. We present additional experiments in the Appendix to show that Director is ill-suited for  
 316 generalization-focused settings: Director still performs well in single environment configurations, but  
 317 its performance deteriorates fast with more training tasks. This indicates poor scalability in terms of  
 318 generalization, a limitation to its application in real-world scenarios.

## 319 5.2 Ablation & Sensitivity Studies

320 In the Appendix, we present ablation results confirming the effectiveness of delusion suppression,  
 321  $k$ -medoids pruning and the effectiveness of spatial abstraction via the local perception field. We also  
 322 provide sensitivity study for the number of checkpoints in each proxy problem.

## 323 5.3 Summary of Experiments

324 Within the scope of the experiments, we conclude that **Skipper** provides benefits for generalization;  
 325 And it achieves better generalization when exposed to more training tasks;

326 From the content presented in the Appendix, we deduce additionally that:

- 327 • Spatial abstraction based on the local perception field is crucial for the scalability of the agents;
- 328 • **Skipper** performs well by reliably decomposing the given tasks, and achieving the sub-tasks  
 329 robustly. Its performance is bottlenecked by the accuracy of the estimated proxy problems as well  
 330 as the checkpoint policies, which correspond to goal generalization and capability generalization,  
 331 respectively, identified in [36]. This matches well with our theory. The proposed delusion suppression  
 332 technique (in Appendix) is effective in suppressing plans with non-existent checkpoints as  
 333 targets, thereby increasing the accuracy of the proxy problems;
- 334 • LEAP fails to generalize well within its original form and can generalize better when combined  
 335 with the ideas proposed in this paper; Director may generalize better only in domains where long  
 336 and informative trajectory collection is possible;
- 337 • We verified empirically that, as expected, **Skipper** is compatible with stochasticity.

## 338 6 Conclusions

339 Building on previous work on spatial abstraction [73], we proposed, analyzed and validated **Skipper**,  
 340 which generalizes its learned skills better than the compared methods, due to its combined spatio-  
 341 temporal abstractions.



342 **References**

- 343 [1] S. Alver and D. Precup. Understanding decision-time vs. background planning in model-based  
344 reinforcement learning. *arXiv preprint arXiv:2206.08442*, 2022.
- 345 [2] S. Alver and D. Precup. Minimal value-equivalent partial models for scalable and robust  
346 planning in lifelong reinforcement learning. In S. Chandar, R. Pascanu, H. Sedghi, and  
347 D. Precup, editors, *Proceedings of The 2nd Conference on Lifelong Learning Agents*, volume  
348 232 of *Proceedings of Machine Learning Research*, pages 548–567. PMLR, 22–25 Aug 2023.
- 349 [3] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin,  
350 O. Pieter Abbeel, and W. Zaremba. Hindsight experience replay. *Advances in neural information  
351 processing systems*, 30, 2017.
- 352 [4] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*,  
353 2016.
- 354 [5] P.-L. Bacon, J. Harb, and D. Precup. The option-critic architecture. In *Proceedings of the AAAI  
355 conference on artificial intelligence*, volume 31, 2017.
- 356 [6] A. Bagaria, J. K. Senthil, and G. Konidaris. Skill discovery for exploration and planning using  
357 deep skill graphs. In *International Conference on Machine Learning*, pages 521–531. PMLR,  
358 2021.
- 359 [7] Y. Bengio. The consciousness prior. *arXiv*, 1709.08568, 2017. [http://arxiv.org/abs/  
360 1709.08568](http://arxiv.org/abs/1709.08568).
- 361 [8] Y. Bengio, N. Léonard, and A. Courville. Estimating or propagating gradients through stochastic  
362 neurons for conditional computation. *arXiv preprint:1308.3432*, 2013.
- 363 [9] M. Chevalier-Boisvert, D. Bahdanau, S. Lahlou, L. Willems, C. Saharia, T. H. Nguyen, and  
364 Y. Bengio. Babyai: A platform to study the sample efficiency of grounded language learning.  
365 *International Conference on Learning Representations*, 2018. [http://arxiv.org/abs/1810.  
366 08272](http://arxiv.org/abs/1810.08272).
- 367 [10] M. Chevalier-Boisvert, L. Willems, and S. Pal. Minimalistic gridworld environment for openai  
368 gym. *GitHub repository*, 2018. <https://github.com/maximecb/gym-minigrid>.
- 369 [11] K. Cobbe, C. Hesse, J. Hilton, and J. Schulman. Leveraging procedural generation to benchmark  
370 reinforcement learning. In *International conference on machine learning*, pages 2048–2056.  
371 PMLR, 2020.
- 372 [12] C. Colas, O. Sigaud, and P.-Y. Oudeyer. How many random seeds? statistical power analysis in  
373 deep reinforcement learning experiments, 2018.
- 374 [13] K. Czechowski, T. Odrzygóźdź, M. Zbysiński, M. Zawalski, K. Olejnik, Y. Wu, Łukasz Kuciński,  
375 and P. Miłoś. Subgoal search for complex reasoning tasks. *arXiv preprint:2108.11204*, 2021.
- 376 [14] W. Dabney, M. Rowland, M. Bellemare, and R. Munos. Distributional reinforcement learning  
377 with quantile regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
378 volume 32, 2018.
- 379 [15] S. Dehaene, H. Lau, and S. Kouider. What is consciousness, and could machines have it?  
380 *Science*, 358, 2020.
- 381 [16] T. G. Dietterich. Hierarchical reinforcement learning with the maxq value function decomposi-  
382 tion. *Journal of artificial intelligence research*, 13:227–303, 2000.
- 383 [17] A. Erraqabi, M. Zhao, M. C. Machado, Y. Bengio, S. Sukhbaatar, L. Denoyer, and A. Lazaric.  
384 Exploration-driven representation learning in reinforcement learning. In *ICML 2021 Workshop  
385 on Unsupervised Reinforcement Learning*, 2021.
- 386 [18] B. Eysenbach, R. R. Salakhutdinov, and S. Levine. Search on the replay buffer: Bridging  
387 planning and reinforcement learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-  
388 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,  
389 volume 32. Curran Associates, Inc., 2019.
- 390 [19] C. Florensa, D. Held, X. Geng, and P. Abbeel. Automatic goal generation for reinforcement  
391 learning agents. In *International conference on machine learning*, pages 1515–1528. PMLR,  
392 2018.

- 393 [20] X. Fu, G. Yang, P. Agrawal, and T. Jaakkola. Learning task informed abstractions. In *International Conference on Machine Learning*, pages 3480–3491. PMLR, 2021.  
394
- 395 [21] S. Fujimoto, H. van Hoof, and D. Meger. Addressing function approximation error in actor-critic  
396 methods, 2018.
- 397 [22] A. Gupta, G. Dar, S. Goodman, D. Ciprut, and J. Berant. Memory-efficient transformers via  
398 top- $k$  attention. *arXiv preprint:2106.06899*, 2021.
- 399 [23] D. Hafner, K.-H. Lee, I. Fischer, and P. Abbeel. Deep hierarchical planning from pixels. In A. H.  
400 Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing  
401 Systems*, 2022.
- 402 [24] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent  
403 dynamics for planning from pixels. In *International conference on machine learning*, pages  
404 2555–2565. PMLR, 2019.
- 405 [25] D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap. Mastering diverse domains through world  
406 models. *arXiv preprint:2301.04104*, 2023.
- 407 [26] C. Hogg, U. Kuter, and H. Muñoz Avila. Learning hierarchical task networks for nondeterministic  
408 planning domains. In *Proceedings of the 21st International Joint Conference on Artificial  
409 Intelligence, IJCAI’09*, page 1708–1714, San Francisco, CA, USA, 2009. Morgan Kaufmann  
410 Publishers Inc.
- 411 [27] D. Y.-T. Hui, M. Chevalier-Boisvert, D. Bahdanau, and Y. Bengio. Babyai 1.1, 2020.
- 412 [28] M. Igl, K. Ciosek, Y. Li, S. Tschitschek, C. Zhang, S. Devlin, and K. Hofmann. Generalization  
413 in reinforcement learning with selective noise injection and information bottleneck. *Advances  
414 in neural information processing systems*, 32, 2019.
- 415 [29] D. Kahneman. *Thinking, fast and slow*. 2017.
- 416 [30] L. Kaufman and P. Rousseeuw. *Partitioning Around Medoids (Program PAM)*, chapter 2, pages  
417 68–125. John Wiley & Sons, Ltd, 1990.
- 418 [31] J. Kim, Y. Seo, and J. Shin. Landmark-guided subgoal generation in hierarchical reinforcement  
419 learning. *arXiv preprint:2110.13625*, 2021.
- 420 [32] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint:1412.6980*,  
421 2014.
- 422 [33] C. A. Knoblock. Automatically generating abstractions for planning. *Artificial intelligence*,  
423 68(2):243–302, 1994.
- 424 [34] G. D. Konidaris and A. G. Barto. Efficient skill learning using abstraction selection. In *IJCAI*,  
425 volume 9, pages 1107–1112, 2009.
- 426 [35] T. D. Kulkarni, K. Narasimhan, A. Saeedi, and J. Tenenbaum. Hierarchical deep reinforcement  
427 learning: Integrating temporal abstraction and intrinsic motivation. *Advances in neural  
428 information processing systems*, 29, 2016.
- 429 [36] L. L. D. Langosco, J. Koch, L. D. Sharkey, J. Pfau, and D. Krueger. Goal misgeneralization in  
430 deep reinforcement learning. In *International Conference on Machine Learning*, volume 162,  
431 pages 12004–12019, 17-23 Jul 2022.
- 432 [37] L. Li, T. J. Walsh, and M. L. Littman. Towards a unified theory of state abstraction for mdps.  
433 *AI&M*, 1(2):3, 2006.
- 434 [38] A. Manchin, E. Abbasnejad, and A. Van Den Hengel. Reinforcement learning with attention  
435 that works: A self-supervised approach. In *Neural Information Processing: 26th International  
436 Conference, ICONIP 2019, Sydney, NSW, Australia, December 12–15, 2019, Proceedings, Part  
437 V 26*, pages 223–230. Springer, 2019.
- 438 [39] A. McGovern and A. G. Barto. Automatic discovery of subgoals in reinforcement learning  
439 using diverse density. 2001.
- 440 [40] I. Momennejad, E. M. Russek, J. H. Cheong, M. M. Botvinick, N. D. Daw, and S. J. Gershman.  
441 The successor representation in human reinforcement learning. *Nature human behaviour*,  
442 1(9):680–692, 2017.

- 443 [41] A. Mott, D. Zoran, M. Chrzanowski, D. Wierstra, and D. Jimenez Rezende. Towards in-  
444 terpretable reinforcement learning using attention augmented agents. *Advances in neural*  
445 *information processing systems*, 32, 2019.
- 446 [42] O. Nachum, S. S. Gu, H. Lee, and S. Levine. Data-efficient hierarchical reinforcement learning.  
447 *Advances in neural information processing systems*, 31, 2018.
- 448 [43] A. Nair, V. Pong, M. Dalal, S. Bahl, S. Lin, and S. Levine. Visual reinforcement learning with  
449 imagined goals. *arXiv preprint:1807.04742*, 2018.
- 450 [44] S. Nair and C. Finn. Hierarchical foresight: Self-supervised learning of long-horizon tasks via  
451 visual subgoal generation. In *International Conference on Learning Representations*, 2020.
- 452 [45] S. Nasiriany, V. Pong, S. Lin, and S. Levine. Planning with goal-conditioned policies. *Advances*  
453 *in Neural Information Processing Systems*, 32, 2019.
- 454 [46] G. Paolo, J. Gonzalez-Billandon, A. Thomas, and B. Kégl. Guided safe shooting: model based  
455 reinforcement learning with safety constraints, 2022.
- 456 [47] A. Patterson, S. Neumann, M. White, and A. White. Empirical design in reinforcement learning,  
457 2023.
- 458 [48] V. Pong, S. Gu, M. Dalal, and S. Levine. Temporal difference models: Model-free deep rl for  
459 model-based control. *arXiv preprint:1802.09081*, 2018.
- 460 [49] M. L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John  
461 Wiley & Sons, 2014.
- 462 [50] R. Y. Rubinstein. Optimization of computer simulation models with rare events. *European*  
463 *Journal of Operational Research*, 99(1):89–112, 1997.
- 464 [51] R. Y. Rubinstein. Optimization of computer simulation models with rare events. *European*  
465 *Journal of Operational Research*, 99(1):89–112, 1997.
- 466 [52] E. D. Sacerdoti. Planning in a hierarchy of abstraction spaces. *Artificial intelligence*, 5(2):115–  
467 135, 1974.
- 468 [53] N. Savinov, A. Dosovitskiy, and V. Koltun. Semi-parametric topological memory for navigation.  
469 *arXiv preprint arXiv:1803.00653*, 2018.
- 470 [54] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lock-  
471 hart, D. Hassabis, T. Graepel, et al. Mastering atari, go, chess and shogi by planning with a  
472 learned model. *Nature*, 588(7839):604–609, 2020.
- 473 [55] D. Shah, P. Xu, Y. Lu, T. Xiao, A. Toshev, S. Levine, and B. Ichter. Value function spaces:  
474 Skill-centric state abstractions for long-horizon reasoning. *arXiv preprint arXiv:2111.03189*,  
475 2021.
- 476 [56] A. Sharma, S. Gu, S. Levine, V. Kumar, and K. Hausman. Dynamics-aware unsupervised  
477 discovery of skills. *arXiv preprint arXiv:1907.01657*, 2019.
- 478 [57] L. X. Shi, J. J. Lim, and Y. Lee. Skill-based model-based reinforcement learning. *arXiv preprint*  
479 *arXiv:2207.07560*, 2022.
- 480 [58] D. Silver and K. Ciosek. Compositional planning using optimal option models. *arXiv preprint*  
481 *arXiv:1206.6473*, 2012.
- 482 [59] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker,  
483 M. Lai, A. Bolton, et al. Mastering the game of go without human knowledge. *Nature*,  
484 550(7676):354–359, 2017.
- 485 [60] K. Sohn, H. Lee, and X. Yan. Learning structured output representation using deep conditional  
486 generative models. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors,  
487 *Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- 488 [61] R. S. Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *SIGART*  
489 *Bulletin*, 2(4):160–163, 1991.
- 490 [62] R. S. Sutton, M. C. Machado, G. Z. Holland, D. S. F. Timbers, B. Tanner, and A. White. Reward-  
491 respecting subtasks for model-based reinforcement learning. *arXiv preprint:2202.03466*, 2022.
- 492 [63] R. S. Sutton, D. Precup, and S. Singh. Between mdps and semi-mdps: A framework for temporal  
493 abstraction in reinforcement learning. *Artificial Intelligence*, 112(1):181–211, 1999.

- 494 [64] T. Sylvain, L. Petrini, and D. Hjelm. Locality and compositionality in zero-shot learning. *arXiv*  
495 *preprint arXiv:1912.12179*, 2019.
- 496 [65] C. Tang and R. R. Salakhutdinov. Multiple futures prediction. *Advances in neural information*  
497 *processing systems*, 32, 2019.
- 498 [66] Y. Tang, D. Nguyen, and D. Ha. Neuroevolution of self-interpretable agents. In *Proceedings of*  
499 *the 2020 Genetic and Evolutionary Computation Conference*, pages 414–424, 2020.
- 500 [67] A. Van Den Oord, O. Vinyals, et al. Neural discrete representation learning. *Advances in neural*  
501 *information processing systems*, 30, 2017.
- 502 [68] H. Van Hasselt, A. Guez, and D. Silver. Deep reinforcement learning with double q-learning.  
503 In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
- 504 [69] K. Xie, H. Bharadhwaj, D. Hafner, A. Garg, and F. Shkurti. Latent skill planning for exploration  
505 and transfer. *arXiv preprint arXiv:2011.13897*, 2020.
- 506 [70] A. Zadaianchuk, M. Seitzer, and G. Martius. Self-supervised visual reinforcement learning with  
507 object-centric representations. *arXiv preprint arXiv:2011.14381*, 2020.
- 508 [71] A. Zhang, C. Lyle, S. Sodhani, A. Filos, M. Kwiatkowska, J. Pineau, Y. Gal, and D. Precup.  
509 Invariant causal prediction for block mdps. In *International Conference on Machine Learning*,  
510 pages 11214–11224. PMLR, 2020.
- 511 [72] L. Zhang, G. Yang, and B. C. Stadie. World model as a graph: Learning latent landmarks  
512 for planning. In *International Conference on Machine Learning*, pages 12611–12620. PMLR,  
513 2021.
- 514 [73] M. Zhao, Z. Liu, S. Luan, S. Zhang, D. Precup, and Y. Bengio. A consciousness-inspired  
515 planning agent for model-based reinforcement learning. In M. Ranzato, A. Beygelzimer,  
516 Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing*  
517 *Systems*, volume 34, pages 1569–1581. Curran Associates, Inc., 2021.
- 518 [74] M. Zhao, S. Luan, I. Porada, X.-W. Chang, and D. Precup. Meta-learning state-based eligibility  
519 traces for more sample-efficient policy evaluation. *arXiv preprint arXiv:1904.11439*, 2019.

## 520 A APPENDIX

521 Please use the following to quickly navigate to your points of interest.

- 522 • **Weaknesses & Limitations** (Sec. B)
- 523 • **Skipper Algorithmic Details** (Sec. C): pseudocodes,  $k$ -medoids based pruning, delusion suppression
- 524
- 525 • **Theoretical Analyses** (Sec. D.1): detailed proofs, discussions
- 526 • **Implementation Details** (Sec. E): for **Skipper**, **LEAP** and **Director**
- 527 • **More Experiments** (Sec. F): experimental results that cannot be presented in the main paper due
- 528 to page limit
- 529 • **Ablation Tests & Sensitivity Analyses** (Sec. G)

## 530 B Weaknesses & Limitations

531 We would like to expand the discussions on the limitations to the current form of **Skipper**, as well as

532 the design choices that we seek to improve in the future:

- 533 • We generate future checkpoints at random by sampling the partial description space. Despite the
- 534 post-processing such as pruning, the generated checkpoints do not prioritize on the predictable,
- 535 important states that matter the most to form a meaningful long-term plan.
- 536 • The current implementation is intended for pixel input fully-observable tasks with discrete state and
- 537 action spaces. Such a minimalistic form is because we wish to isolate the unwanted challenges from
- 538 other factors that are not closely related to the idea of this work, as well as to make the agent as
- 539 generalist as possible. **Skipper** is naturally compatible with continuous actions spaces and the only
- 540 thing we will need to do is to replace the baseline agent with a compatible one such as TD3 [21];
- 541 on the other hand, for continuous state spaces, the identification of the achievement of a checkpoint
- 542 becomes tricky. This is due to the fact that a strict identity between the current state and the target
- 543 checkpoint may be ever established, we either must adopt a distance measure for approximate state
- 544 equivalence, or rely on the equivalence of the partial descriptions (which is adopted in the current
- 545 implementation). We intentionally designed the partial descriptions to be in the form of bundles
- 546 of binary variables, so that this comparison could be done fast and trivially for any forms of the
- 547 state space; for partial observability, despite that no recurrent mechanism has been incorporated
- 548 in the current implementation, the framework is not incompatible. To implement that, we will
- 549 need to augment the state encoder with recurrent or memory mechanisms and we need to make
- 550 the checkpoint generator directly work over the learned state representations. We acknowledge
- 551 that future work is needed to verify **Skipper**'s performance on the popular partially-observable
- 552 benchmark suites, which requires the incorporation of components to handle partial observability
- 553 as well as scaling up the architectures for more expressive power;
- 554 • We do not know the precise boundaries of the motivating theory on proxy problems, since it only
- 555 indicates performance guarantees on the condition of estimation accuracy, which in turn does not
- 556 correspond trivially to a set of well-defined problems. We are eager to explore, outside the scope of
- 557 sparse-reward navigation, how this approach can be used to facilitate better generalization, and at
- 558 the same time, try to find more powerful theories that guide us better;

## 559 C Skipper's Algorithmic Details

### 560 C.1 Overall Skipper Framework (Pseudo-Code)

561 The pseudocode of **Skipper** is provided in Alg. 1, together with the hyperparameters used in our

562 implementation.

### 563 C.2 $k$ -medoids based pruning

564 We present the pseudocode of the modified  $k$ -medoids algorithm for pruning overcrowded checkpoints

565 in Alg. 2. Note that the presented pseudocode is optimized for readers' understanding, while the

---

**Algorithm 1: Skipper** with Random Checkpoints (implementation choice in purple)

---

```
for each episode do
  // — start of the subroutine to construct the proxy problem
  generate more than necessary (32) checkpoints by sampling from the partial descriptions
  given the extracted context from the initial state;
  ( $k = 12$ )-medoid pruning upon estimated distances among all checkpoints; // prune vertices
  use estimators to annotate the edges between the nodes (including a terminal state estimator
  to correct the estimates);
  prune edges that are too far-fetched according to distance estimations (threshold set to be 8,
  same as replan interval); // prune edges
  // — end of the subroutine to construct the proxy problem
  for each agent-environment interaction step until termination of episode do
    if decided to explore (DQN-style annealing  $\epsilon$ -greedy) then
      take a random action;
    else
      if abstract problem just constructed or a checkpoint / timeout reached ( $\geq 8$  steps
      since last planned) then
        [OPTIONAL] call the subroutine above for Skipper-regen;
        run value iteration (for 5 iterations) on the proxy problem, select the target
        checkpoint;
      follow the action suggested by the checkpoint-achieving policy;
      if time to train (every 4 actions) then
        sample hindsight transitions and train checkpoint-achieving policy, estimators
        (including a terminal state estimator) and checkpoint generator;
        [OPTIONAL]: train estimators with generated checkpoints to suppress delusion;
      save interaction into the trajectory experience replay;
    convert trajectory into HER samples (relabel 4 random states as additional goals);
```

---

566 actual implementation is parallelized. The changes upon the original  $k$ -medoids algorithm is marked  
567 in purple, which implement a forced preservation of data points: when  $k$ -medoids is called after the  
568 unpruned graph is constructed,  $\mathcal{S}_v$  is set to be the set containing the goal state only. This is intended  
569 to span more uniformly in the state space with checkpoints, while preserving the goal.

570 Let the estimated distance matrix be  $D$ , where each element  $d_{i,j}$  represents the estimated trajectory  
571 length it takes for  $\pi$  to fulfill the transition from checkpoint  $i$  to checkpoint  $j$ . Since  $k$ -medoids  
572 cannot handle infinite distances (e.g. from a terminal state to another state), the distance matrix  $D$  is  
573 truncated, and then we take the elementwise minimum between the truncated  $D$  and  $D^T$  to preserve  
574 the one-way distances. The matrix containing the elementwise minimums would be the input of the  
575 pruning algorithm.

### 576 C.3 Delusion Suppression

577 RL agents are prone to blindly optimizing for an intrinsic objective without fully understanding the  
578 consequences of its actions. Particularly in model-based RL or in Hierarchical RL (HRL), there is  
579 a significant risk posed by the agents trying to achieve delusional future states that do not exist or  
580 beyond the safety constraints. With a use of a learned generative model, as in **Skipper** and other HP  
581 frameworks, such risk is almost inevitable, because of uncontrollable generalization effects.

582 Generalization abilities of the generative models are a double-edged sword. The agent would take  
583 advantage of its potentials to propose novel checkpoints to improve its behavior, but is also at risk  
584 of wanting to achieve non-existent unknown consequences. In **Skipper**, checkpoints imagined by  
585 the generative model could correspond to non-existent “states” that would lead to delusional edge  
586 estimates and therefore confuse planning. For instance, arbitrarily sampling partial descriptions may  
587 result in a delusional state where the agent is in a cell that can never be reached from the initial states.  
588 Since such states do not exist in the experience replay, the estimators will have not learned how to  
589 handle them appropriately when encountered in the generated proxy problem during decision time.  
590 We present a resulting failure mode in Fig. 4.

---

**Algorithm 2:** Checkpoint Pruning with  $k$ -medoids

---

**Data:**  $X = \{x_1, x_2, \dots, x_n\}$  (state indices),  $D$  (estimated distance matrix),  $\mathcal{S}_v$  (states that must be kept),  $k$  (#checkpoints to keep)

**Result:**  $\mathcal{S}_\odot \equiv \{M_1, M_2, \dots, M_k\}$  (checkpoints kept)

Initialize  $\mathcal{S}_\odot \equiv \{M_1, M_2, \dots, M_k\}$  randomly from  $X$

make sure  $\mathcal{S}_v \subset \mathcal{S}_\odot$

repeat

    Assign each data point  $x_i$  to the nearest medoid  $M_j$ , forming clusters  $C_1, C_2, \dots, C_k$ ;

**foreach** medoid  $M_j$  **do**

        Calculate the cost  $J_j$  of  $M_j$  as the sum of distances between  $M_j$  and the data points in  $C_j$ ;

    Find the medoid  $M_j$  with the lowest cost  $J_j$ ;

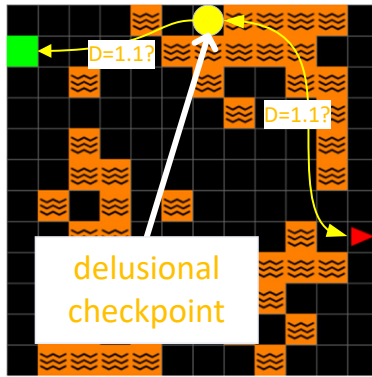
**if**  $M_j$  changes **then**

        make sure  $\mathcal{S}_v \subset \mathcal{S}_\odot$

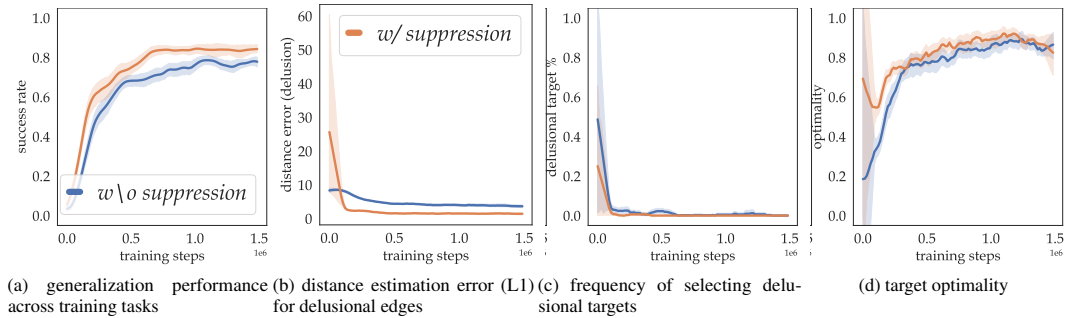
        Replace  $M_j$  with the data point in  $C_j$  that minimizes the total cost;

**until** Convergence (no cost improvement);

---



**Figure 4: Example of Failure Caused by Delusions:** we illustrate an instance of chasing delusional checkpoint in one of our experimental runs by **Skipper**. The distance (discount) estimator, probably due to the ill-generalization, estimates that the delusional checkpoint (yellow) is very close to every other state. A resulting plan was that the agent thought it could reach any far-away checkpoints by using the delusional state to form a shortcut: the goal that was at least 17 steps away would be reached in 2.2.



**Figure 5: Performance of Skipper-once with the proposed Delusion Suppression Technique:** each curve and corresponding error bar (95% CI) are processed from 20 independent seed runs. a) the performance across training tasks is shown. A more optimal performance can be achieved with **Skipper**-once in training tasks, when delusions are suppressed; b) During training interactions, the error in estimated (truncated) distance from and to delusional targets are significantly reduced with the technique; c) The frequency of selecting a delusional target is reduced to almost negligible during the whole training process; d) The optimality of target checkpoint during training can be improved by the suppression. Each agent is trained with 50 environments and each curve is processed from 20 independent seed runs.

591 To address such concerns, we propose an optional auxiliary training procedure that makes the agent  
 592 stay further away from delusional checkpoints. Due to the favorable properties of the update rules  
 593 of  $D_\pi$  (in fact,  $V_\pi$  as well), all we have to do is to replace the hindsight-sampled target states  
 594 with generated checkpoints, which contain non-existent states. Then, the auxiliary rewards will all  
 595 converge to the minimum in terms of favorability on the non-existent states. This is implemented  
 596 trivially by adding a loss to the original training loss for the distance estimator, which we give a 0.25  
 597 scaling for stability.

---

**Algorithm 3:** Delusion Suppression

---

// This whole code block should be injected into the training loop if used  
 generate using the checkpoint generator, from the sampled batch of encoded states, the target  
 states (to overwrite those relabelled in the HER) *i.e.* replace  $\langle s_t, a_t, r_{t+1}, s_{t+1}, s^\odot \rangle$  with  
 $\langle s_t, a_t, r_{t+1}, s_{t+1}, s_*^\odot \rangle$ , where  $s_*^\odot$  are generated from the context of  $s_t$   
 train the distance estimator  $D$  as if these are sampled from the HER

---

598 We provide analytic results and related discussion for **Skipper-once** agents trained with the proposed  
 599 delusion suppression technique on 50 training tasks in Fig. 5. The delusion suppression technique is  
 600 not enabled by default because it was not introduced in the main manuscript due to the page limits.

601 The delusion suppression technique can also be used to help us understand the failure modes of LEAP  
 602 in Sec. E.2.4.

## 603 D Theoretical Analyses

### 604 D.1 Update Rules for Edge Estimation

605 First, we want to show that the update rules proposed in the main paper indeed estimate the desired  
 606 cumulative discount and reward.

607 The low-level checkpoint-achieving policy  $\pi$  is trained with an intrinsic reward to reach target state  
 608  $s^\odot$ . The cumulative reward and cumulative discount are estimated by applying policy evaluation  
 609 given  $\pi$ , on the two sets of auxiliary reward signals, respectively.

610 For the cumulative discounted reward random variable:

$$V_\pi(s_t, a_t | s^\odot) = R(s_t, a_t, S_{t+1}) + \gamma V_\pi(S_{t+1}, A_{t+1} | s^\odot) \quad (5)$$

$$= \sum_{\tau=t}^{\infty} \gamma^{\tau-t} R(S_\tau, A_\tau, S_{\tau+1}), \quad (6)$$

611 where  $S_{t+1} \sim p(\cdot | s_t, a_t)$ ,  $A_{t+1} \sim \pi(\cdot | S_{t+1}, s^\odot)$ , and with  $V_\pi(S_{t+1}, A_{t+1} | s^\odot) = 0$  if  $S_{t+1} = s^\odot$ .  
 612 We overload the notation as follows:  $V_\pi(s | s^\odot) \doteq V_\pi(s, A | s^\odot)$  with  $A \sim \pi(\cdot | s, s^\odot)$ .

613 The cumulative discount random variable denotes the event that the trajectory did not terminate  
 614 before reaching the target  $s^\odot$ :

$$\Gamma_\pi(S_t, A_t | s^\odot) = \gamma \cdot \Gamma_\pi(S_{t+1}, A_{t+1} | s^\odot), \quad (7)$$

$$= \gamma^{T_\perp - t} \mathbb{I}\{S_{T_\perp} = s^\odot\}, \quad (8)$$

615 where  $T_\perp$  denotes the timestep when the trajectory terminates, and with  $\Gamma_\pi(S_{t+1}, A_{t+1} | s^\odot) = 1$   
 616 if  $S_{t+1} = s^\odot$  and  $\Gamma_\pi(S_{t+1}, A_{t+1} | s^\odot) = 0$  if  $S_{t+1} \neq s^\odot$  is terminal. We overload the notation as  
 617 follows:  $\Gamma_\pi(s_t | s^\odot) \doteq \Gamma_\pi(s_t, A_t | s^\odot)$  with  $A_{t+1} \sim \pi(\cdot | S_{t+1}, s^\odot)$ .

618 Note that, for the sake of simplicity, we take here the view that the terminality of states is deterministic,  
 619 but this is not reductive as any state with a stochastic terminality can be split into two identical states:  
 620 one that is deterministically non-terminal and the other that is deterministically terminal. Note also  
 621 that we could adopt the view that the discount factor is the constant probability of the trajectory to  
 622 not terminate.

### 623 D.2 Performance Bound

624 We are going to denote the expected cumulative discounted reward, *a.k.a.* the state-action value  
 625 with  $q_\pi \doteq \mathbb{E}_\pi[V]$ , and let  $\hat{q}_\pi$  be our estimate for it. We are also going to consider the state value



626  $v_\pi(s|s^\odot) \doteq \sum_a \pi(a|s, s^\odot) q_\pi(s, a|s^\odot)$  and its estimate  $\hat{v}_\pi$ . Similarly, we denote the expected  
627 cumulative discount with  $\gamma_\pi \doteq \mathbb{E}_\pi[\Gamma]$  and its estimate with  $\hat{\gamma}_\pi$ .

628 We are in the presence of a hierarchical policy. The high level policy  $\mu$  consists in (potentially)  
629 stochastically picking a sequence of checkpoints. The low-level policy is implemented by  $\pi$  which  
630 is assumed to be given and fixed for the moment. The composite policy  $\mu \circ \pi$  is non-Markovian: it  
631 depends both on the current state and the current checkpoint goal. So there is no notion of state value,  
632 except when we arrive at a checkpoint, *i.e.* when a high level action (checkpoint selection) needs to  
633 be chosen.

634 Proceeding further, we adopt the view where the discounts are a way to represent the hazard of the  
635 environment:  $1 - \gamma$  is the probability of sudden trajectory termination. In this view,  $v_\pi$  denotes  
636 the (undiscounted: there is no more discounting) expected sum of reward before reaching the next  
637 checkpoint, and more interestingly  $\gamma_\pi$  denotes the binomial random variable of non-termination  
638 during the transition to the selected checkpoint.

639 Making the following assumption that the trajectory terminates almost surely when reaching the goal,  
640 *i.e.*  $\gamma_\pi(s_i, s_g) = 0, \forall s_i$ , the gain  $V$  can be written:

$$V_0 = V(S_0^\odot | S_1^\odot) + \Gamma(S_0^\odot | S_1^\odot) V_1 = \sum_{k=0}^{\infty} V(S_k^\odot | S_{k+1}^\odot) \prod_{i=0}^{k-1} \Gamma(S_i^\odot | S_{i+1}^\odot), \quad (9)$$

641 where  $S_{k+1} \sim \mu(\cdot | S_k)$ , where  $V(S_k^\odot | S_{k+1}^\odot)$  is the gain obtained during the path between  $S_k^\odot$  and  
642 where  $S_{k+1}^\odot$ , and  $\Gamma(S_k^\odot | S_{k+1}^\odot)$  is either 0 or 1 depending whether the trajectory terminated or reached  
643  $S_{k+1}^\odot$ . If we consider  $\mu$  as a deterministic planning routine over the checkpoints, then the action space  
644 of  $\mu$  boils down to a list of checkpoints  $\{s_0^\odot = s_0, s_1^\odot, \dots, s_n^\odot = s_g\}$ . Thanks to the Markovian  
645 property in checkpoints, we have independence between  $V_\pi$  and  $\Gamma_\pi$ , therefore for the expected value  
646 of  $\mu \circ \pi$ , we have:

$$v_{\mu \circ \pi}(s_0) \doteq \mathbb{E}_{\mu \circ \pi}[V | S_0 = s_0] = \sum_{k=0}^{\infty} v_\pi(s_k^\odot | s_{k+1}^\odot) \prod_{i=0}^{k-1} \gamma_\pi(s_i^\odot | s_{i+1}^\odot) \quad (10)$$

647 Having obtained the ground truth value, in the following, we are going to consider the estimates  
648 which may have small error terms:

$$|v_\pi(s) - \hat{v}_\pi(s)| < \epsilon_v v_{\max} \ll (1 - \gamma) v_{\max} \quad \text{and} \quad |\gamma_\pi(s) - \hat{\gamma}_\pi(s)| < \epsilon_\gamma \ll (1 - \gamma)^2 \quad \forall s. \quad (11)$$

649 We are looking for a performance bound, and assume without loss of generality that the reward  
650 function is non-negative, *s.t.* the values are guaranteed to be non-negative as well. We provide an  
651 upper bound:

$$\hat{v}_{\mu \circ \pi}(s) \doteq \sum_{k=0}^{\infty} \hat{v}_\pi(s_k^\odot | s_{k+1}^\odot) \prod_{i=0}^{k-1} \hat{\gamma}_\pi(s_i^\odot | s_{i+1}^\odot) \quad (12)$$

$$\leq \sum_{k=0}^{\infty} (v_\pi(s_k^\odot | s_{k+1}^\odot) + \epsilon_v v_{\max}) \prod_{i=0}^{k-1} (\gamma_\pi(s_i^\odot | s_{i+1}^\odot) + \epsilon_\gamma) \quad (13)$$

$$\leq v_{\mu \circ \pi}(s) + \sum_{k=0}^{\infty} \epsilon_v v_{\max} \prod_{i=0}^{k-1} (\gamma_\pi(s_i^\odot | s_{i+1}^\odot) + \epsilon_\gamma) + \sum_{k=0}^{\infty} (v_\pi(s_k^\odot | s_{k+1}^\odot) + \epsilon_v v_{\max}) k \epsilon_\gamma \gamma^k + o(\epsilon_v + \epsilon_\gamma) \quad (14)$$

$$\leq v_{\mu \circ \pi}(s) + \epsilon_v v_{\max} \sum_{k=0}^{\infty} \gamma^k + \epsilon_\gamma v_{\max} \sum_{k=0}^{\infty} k \gamma^k + o(\epsilon_v + \epsilon_\gamma) \quad (15)$$

$$\leq v_{\mu \circ \pi}(s) + \frac{\epsilon_v v_{\max}}{1 - \gamma} + \frac{\epsilon_\gamma v_{\max}}{(1 - \gamma)^2} + o(\epsilon_v + \epsilon_\gamma) \quad (16)$$

652 Similarly, we can derive a lower bound:

$$\hat{v}_{\mu \circ \pi}(s) \doteq \sum_{k=0}^{\infty} \hat{v}_{\pi}(s_k^{\circ} | s_{k+1}^{\circ}) \prod_{i=0}^{k-1} \hat{\gamma}_{\pi}(s_i^{\circ} | s_{i+1}^{\circ}) \quad (17)$$

$$\geq \sum_{k=0}^{\infty} (v_{\pi}(s_k^{\circ} | s_{k+1}^{\circ}) - \epsilon_v v_{\max}) \prod_{i=0}^{k-1} (\gamma_{\pi}(s_i^{\circ} | s_{i+1}^{\circ}) - \epsilon_{\gamma}) \quad (18)$$

$$\geq v_{\mu \circ \pi}(s) - \sum_{k=0}^{\infty} \epsilon_v v_{\max} \prod_{i=0}^{k-1} (\gamma_{\pi}(s_i^{\circ} | s_{i+1}^{\circ}) - \epsilon_{\gamma}) - \sum_{k=0}^{\infty} (v_{\pi}(s_k^{\circ} | s_{k+1}^{\circ}) - \epsilon_v v_{\max}) k \epsilon_{\gamma} \gamma^k + o(\epsilon_v + \epsilon_{\gamma}) \quad (19)$$

$$\geq v_{\mu \circ \pi}(s) - \epsilon_v v_{\max} \sum_{k=0}^{\infty} \gamma^k - \epsilon_{\gamma} v_{\max} \sum_{k=0}^{\infty} k \gamma^k + o(\epsilon_v + \epsilon_{\gamma}) \quad (20)$$

$$\geq v_{\mu \circ \pi}(s) - \frac{\epsilon_v v_{\max}}{1 - \gamma} - \frac{\epsilon_{\gamma} v_{\max}}{(1 - \gamma)^2} + o(\epsilon_v + \epsilon_{\gamma}) \quad (21)$$

653 We may therefore conclude that  $\hat{v}_{\mu \circ \pi}$  equals  $v_{\mu \circ \pi}$  up to an accuracy of  $\frac{\epsilon_v v_{\max}}{1 - \gamma} + \frac{\epsilon_{\gamma} v_{\max}}{(1 - \gamma)^2} + o(\epsilon_v + \epsilon_{\gamma})$ .  
 654 Note that the requirement for the reward function to be positive is only a cheap technical trick to  
 655 ensure we bound in the right direction of  $\epsilon_{\gamma}$  errors in the discounting, but that the theorem would still  
 656 stand if it were not the case.

### 657 D.3 No Assumption on Optimality

658 If the low-level policy  $\pi$  is perfect, then the best high-level policy  $\mu$  is to choose directly the goal  
 659 as target<sup>3</sup>. Our approach assumes that it would be difficult to learn effectively a  $\pi$  when the target  
 660 is too far, and that we would rather use a proxy to construct a path with shorter-distance transitions.  
 661 Therefore, we’ll never want to make any optimality assumption on  $\pi$ , otherwise our approach is  
 662 pointless. These theories we have initiated makes no assumption on  $\pi$ .

663 The Theorem provides guarantees on the solution to the overall problem. The quality of the solution  
 664 depends on both the quality of the estimates (distances/discounts, rewards) and the quality of the  
 665 policy, as the theorem guarantees accuracy to the solution of the overall problem given a current  
 666 policy, which should evolve towards optimal during training. This means bad policy with good  
 667 estimation will lead to an accurate yet bad overall solution. No matter the quality of the policy, with a  
 668 bad estimation, it will result in a poor estimate of solutions. Only a near-optimal policy and good  
 669 estimation will lead to a near-optimal solution.

## 670 E Implementation Details for Experiments

### 671 E.1 Skipper

#### 672 E.1.1 Training

673 The agent is based on a distributional prioritized double DQN. All the trainable parameters are  
 674 optimized with Adam at a rate of  $2.5 \times 10^{-4}$  [32], with a gradient clipping by value (maximum  
 675 absolute value 1.0). The priorities for experience replay sampling are equal to the per-sample training  
 676 loss.

#### 677 E.1.2 Full State Encoder

678 The full-state encoder is a two layered residual block (with kernel size 3 and doubled intermediate  
 679 channels) combined with the 16-dimensional bag-of-words embedder of BabyAI [27].

<sup>3</sup>A triangular inequality can be shown that with a perfect  $\pi$  and a perfect estimate of  $v_{\pi}$  and  $\gamma_{\pi}$ , the performance will always be minimized by selecting  $s_1^{\circ} = s_g$ .

680 **E.1.3 Partial State Selector (Spatial Abstraction)**

681 The selector  $\sigma$  is implemented with one-head (not multiheaded, therefore the output linear transfor-  
682 mation of the default multihead attention implementation in PyTorch is disabled.) top-4 attention,  
683 with each local perceptive field of size  $8 \times 8$  cells. Layer normalization [4] is used before and after  
684 the spatial abstraction.

685 **E.1.4 Estimators**

686 The estimators, which operate on the partial states, are 3-layered  
687 MLPs with 256 hidden units.

688 An additional estimator for termination is learned, which instead  
689 of taking a pair of partial states as input, takes only one, and is  
690 learned to classify terminal states with cross-entropy loss. The  
691 estimated distance from terminal states to other states would be  
692 overwritten with  $\infty$ . The internal  $\gamma$  for intrinsic reward of  $\pi$  is  
693 0.95, while the task  $\gamma$  is 0.99

694 The estimators use C51 distributional TD learning [14]. That is,  
695 the estimators output histograms (softmax over vector outputs)  
696 instead of scalars. We regress the histogram towards the targets,  
697 where these targets are skewed histograms of scalar values, towards  
698 which KL-divergence is used to train. At the output, there are  
699 16 bins for each histogram estimation (value for policy, reward,  
700 distance).

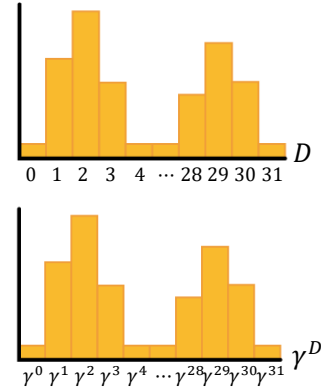


Figure 6: **Estimating Distributions of Discount and Distance with the Same Histogram:** by transplanting the support with the corresponding discount values, the distribution of the cumulative discount can be inferred.

701 **E.1.5 Recovering Discounts from Distances**

702 We recover the distribution of the cumulative discount by replacing  
703 the support of the discretized truncated distances with the corre-  
704 sponding discounts, as illustrated in Fig. 6. This addresses the  
705 problem of  $\mathbb{E}[\gamma^D] \neq \gamma^{\mathbb{E}[D]}$ , as the probability of having a trajec-  
706 tory length of 4 under policy  $\pi$  from state  $s_t$  to  $s_\odot$  is the same as a  
707 trajectory having discount  $\gamma^4$ .

708 **E.1.6 Checkpoint Generator**

709 Despite **Skipp**er is designed to have the generator work on state level, that is, it should take learned  
710 state representations as inputs and have state representations as outputs, in our experiments, the  
711 generator actually operates on observation inputs and outputs. This is because of the preferred  
712 compactness of the observations and the equivalence to full states under full observability in our  
713 experiments.

714 The context extractor  $\mathcal{E}_c$  is a 32-dimensional BabyAI BOW embedder. It encodes an input observation  
715 into a representation of the episodic context.

716 The partial description extractor  $\mathcal{E}_z$  is made of a 32-dimensional BabyAI BOW embedder, followed by  
717 3 aforementioned residual blocks with  $3 \times 3$  convolutions (doubling the feature dimension every time)  
718 in between, ended by global maxpool and a final linear projection to the latent weights. The partial  
719 descriptions are bundles of 6 binary latents, which could represent at most 64 “kinds” of checkpoints.  
720 Inspired by VQ-VAE [67], we use the argmax of the latent weights as partial descriptions, instead of  
721 sampling according to the softmax-ed weights. This enables easy comparison of current state to the  
722 checkpoints in the partial description space, because each state deterministically corresponds to one  
723 partial description. We identify reaching a target checkpoint if the partial description of the current  
724 state matches that of the target.

725 The fusing function first projects linearly the partial descriptions to a 128-dimensional space and then  
726 uses deconvolution to recover an output which shares the same size as the encoded context. Finally,  
727 a residual block is used, followed by a final  $1 \times 1$  convolution that downscales the concatenation of  
728 context together with the deconv’ed partial description into a 2D weight map. The agent’s location is  
729 taken to be the argmax of this weight map.

730 The whole checkpoint generator is trained end-to-end with a standard VAE loss. That is the sum  
731 of a KL-divergence for the agent’s location, and the entropy of partial descriptions, weighted by  
732  $2.5 \times 10^{-4}$ , as suggested in <https://github.com/AntixK/PyTorch-VAE>. Note that the per-  
733 sample losses in the batches are not weighted for training according to priority from the experience  
734 replay.

735 We want to mention that if one does not want to generate non-goal terminal states as checkpoints, we  
736 could also seek to train on reversed  $\langle S^\odot, S_t \rangle$  pairs. In this case, the checkpoints to reconstruct will  
737 never be terminal.

### 738 E.1.7 HER

739 Each experienced transition is further duplicated into 4 hindsight transitions at the end of each episode.  
740 Each of these transitions is combined with a randomly sampled observation from the same trajectory  
741 as the relabelled “goal”. The size of the hindsight buffer is extended to 4 times that of the baseline  
742 that does not learn from hindsight accordingly, that is,  $4 \times 10^6$ .

### 743 E.1.8 Planning

744 As introduced, we use value iteration over options [63] to plan over the proxy problem represented  
745 as an SMDP. We use the matrix form  $Q = R_{S \times S} + \Gamma V$ , where  $R$  and  $\Gamma$  are the estimated edge  
746 matrices for cumulative rewards, respectively. Note that this notation is different from the ones we  
747 used in the manuscript. The checkpoint value  $V$ , initialized as all-zero, is taken on the maximum  
748 of  $Q$  along the checkpoint target (the actions for  $\mu$ ) dimension. When planning is initiated during  
749 decision time, the value iteration step is called 5 times. We do not run until convergence since with  
750 low-quality estimates during the early stages of the learning, this would be a waste of time. The edges  
751 from the current state towards other states are always set to be one-directional, and the self-loops are  
752 also removed. This means the first column as well as the diagonal elements of  $R$  and  $\Gamma$  are all zeros.  
753 Besides pruning edges based on the distance threshold, as introduced in the main paper, the terminal  
754 estimator is also used to prune the matrices  $R$  and  $\Gamma$ : the rows corresponding to the terminal states  
755 are all zeros.

756 The only difference between the two variants, *i.e.* **Skipper-once** and **Skipper-regen** is that the latter  
757 variant would discard the previously constructed proxy problem and construct a new one every time  
758 the planning is triggered. This introduces more computational effort while lowering the chance that  
759 the agent gets “trapped” in a bad proxy problem that cannot form effective plans to achieve the goal.  
760 If such a situation occurs with **Skipper-regen**, as long as the agent does not terminate the episode  
761 prematurely, a new proxy problem will be generated to hopefully address the issue. Empirically, as  
762 we have demonstrated in the experiments, such variant in the planning behavior results in generally  
763 significant improvements in terms of generalization abilities at the cost of extra computation.

### 764 E.1.9 Hyperparameter Tuning

765 Some hyperparameters introduced by **Skipper** can be located in the pseudocode in Alg. 1.

766 **Timeout and Pruning Threshold** Intuitively, we tied the timeout to be equal to the distance pruning  
767 threshold. The timeout kicks in when the agent thinks a checkpoint can be achieved within *e.g.* 8  
768 steps, but already spent 8 steps yet still could not achieve it.

769 This leads to how we tuned the pruning (distance) threshold: we fully used the advantage of our  
770 experiments on DP-solvable tasks: with a snapshot of the agent during its training, we can sample  
771 many  $\langle \text{starting state, target state} \rangle$  pairs and calculate the ground truth distance between the pair, as  
772 well as the failure rate of reaching from the starting state to the target state given the current policy  $\pi$ ,  
773 then plot them as the  $x$  and  $y$  values respectively for visualization. We found such curves to evolve  
774 from high failure rate at the beginning, to a monotonically increasing curve, where at small true  
775 distances, the failure rates are near zero. We picked 8 because the curve starts to grow explosively  
776 when the true distances are more than 9.

777  **$k$  for  $k$ -medoids** We tuned this by running a sensitivity analysis on **Skipper** agents with different  $k$ ’s,  
778 whose results are presented previously in this Appendix.

779 Additionally, we prune from 32 checkpoints because 32 checkpoints could achieve (visually) a good  
780 coverage of the state space as well as its friendliness to NVIDIA accelerators.

781 **Size of local Perception Field** We used a local perception field of size 8 because our baseline model-  
782 free agent would be able to solve and generalize well within  $8 \times 8$  tasks, but not larger. Roughly  
783 speaking, our spatial abstraction breaks down the overall tasks into  $8 \times 8$  sub-tasks, which the policy  
784 could comfortably solve.

785 **Model-free Baseline Architecture** The baseline architecture (distributional, Double DQN) was  
786 heavily influenced by the architecture used in the previous work [73], which demonstrated success  
787 on similar but smaller-scale experiments ( $8 \times 8$ ). The difference is that while then we used compu-  
788 tationally heavy components such as transformer layers on a set-based representation, we replaced  
789 them with a simpler and effective local perception component. We validated our model-free baseline  
790 performance on the tasks proposed in [73].

## 791 E.2 LEAP

### 792 E.2.1 Adaptation for Discrete Action Spaces

793 The LEAP baseline has been implemented from scratch for our experiments, since the original  
794 open-sourced implementation<sup>4</sup> was not compatible with environments with discrete action spaces.  
795 LEAP’s training involves two pretraining stages, that are, generator pretraining and distance estimator  
796 pretraining, which were originally named the VAE and RL pretrainings. Despite our best effort, that  
797 is to be covered in detail, we found that LEAP was unable to get a reasonable performance in its  
798 original form after rebasing it on a discrete model-free RL baseline.

### 799 E.2.2 Replacing the Model

800 We tried to identify the reasons why the generalization performance of the adapted LEAP was  
801 unsatisfactory: we found that the original VAE used in LEAP is not capable to handle even few  
802 training tasks, let alone generalize well to the evaluation tasks. Even by combining the idea of the  
803 context / partial description split (still with continuous latents), during decision time, the planning  
804 results given by the evolutionary algorithm (Cross Entropy Method, CEM, [51]) almost always  
805 produce delusional plans that are catastrophic in terms of performance. This was why we switched  
806 into LEAP the same conditional generator we proposed in the paper, and adapted CEM accordingly,  
807 due to the change from continuous latents to discrete.

808 We also did not find that using the pretrained VAE representation as the state representation during  
809 the second stage helped the agent’s performance, as the paper claimed. In fact, the adapted LEAP  
810 variant could only achieve decent performance after learning a state representation from scratch in  
811 the RL pretraining phase. Adopting **Skipper**’s splitting generator also disables such choice.

### 812 E.2.3 Replacing TDM

813 The original distance estimator based on Temporal Difference Models (TDM) also does not show  
814 capable performance in estimating the length of trajectories, even with the help of a ground truth  
815 distance function (calculated with DP). Therefore, we switched to learning the distance estimates  
816 with our proposed method. Our distance estimator is not sensitive to the sub-goal time budget as  
817 TDM and is hence more versatile in environments like that was used in the main paper, where the  
818 trajectory length of each checkpoint transition could highly vary. Like for **Skipper**, an additional  
819 terminal estimator has been learned to make LEAP planning compatible with the terminal lava states.  
820 Note that this LEAP variant was trained on the same sampling scheme with HER as in **Skipper**.

821 The introduced distance estimator, as well as the accompanying full-state encoder, are of the same  
822 architecture, hyperparameters, and training method as those used in **Skipper**. The number of  
823 intermediate subgoals for LEAP planning is tuned to be 3, which close to how many intermediate  
824 checkpoints **Skipper** typically needs to reach before finishing the tasks. The CEM is called with 5  
825 iterations for each plan construction, with a population size of 128 and an elite population of size  
826 16. We found no significant improvement in enlarging the search budget other than additional wall  
827 time. The new initialization of the new population is by sampling a  $\epsilon$ -mean of the elite population  
828 (the binary partial descriptions), where  $\epsilon = 0.01$  to prevent the loss of diversity. Because of the  
829 very expensive cost of using CEM at decision time and its low return of investment in terms of

---

<sup>4</sup><https://github.com/snasiriany/leap>

Table 1: The changed parameters and their values in the config file of the Director agent.

Parameter	Value
replay_size	2M
replay_chunk	12
imag_horizon	8
env_skill_duration	4
train_skill_duration	4
worker_rews	{extr: 0.5, expl: 0.0, goal: 1.0}
sticky	False
gray	False

830 generalization performance, during the RL pretraining phase, the agent performs random walks over  
 831 uniformly random initial states to collect experience.

### 832 E.2.4 Failure Mode: Delusional Plans

833 Interestingly, we find that a major reason why LEAP does not generalize well is that it often  
 834 generates delusional plans that lead to catastrophic subgoal transitions. This is likely because of its  
 835 blind optimization in the latent space towards shorter path plans: any paths with delusional shorter  
 836 distances would be preferred. We present the results with LEAP combined with our proposed delusion  
 837 suppression technique in Fig. 7. We find that the adapted LEAP agent, with our generator, our  
 838 distance estimator, and the delusion suppression technique, is actually able to achieve significantly  
 839 better generalization performance.

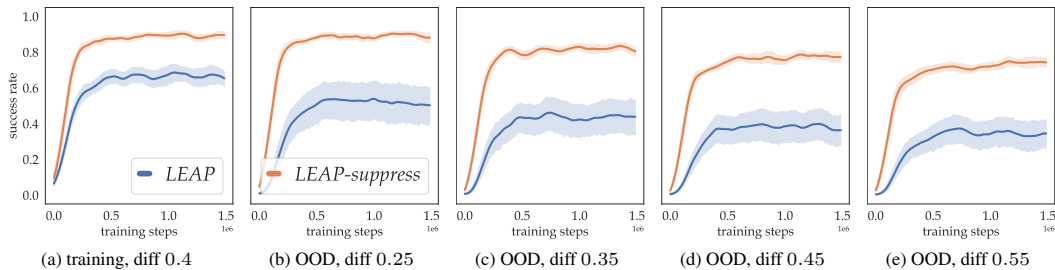


Figure 7: **Comparative Results of LEAP with and without the delusion suppression technique:** the results are obtained with 50 training tasks. The results are obtained from 20 independent seed runs.

## 840 E.3 Director

### 841 E.3.1 Adaptation

842 We based our experiments of Director [23] on the publicly available code  
 843 (<https://github.com/danijar/director>) released by the authors.  
 844 Except for a few changes in the parameters, which are depicted in Tab. 1,  
 845 we have used the default configuration provided for Atari environments.  
 846 Note that as the Director version in which the worker receives no task  
 847 rewards performed worse in our tasks, we have used the version in which  
 848 the worker receives scaled task rewards (referred to as “Director (worker  
 849 task reward)” in [23]). This agent has also been shown to perform better  
 850 across various domains in [23].

851 **Encoder.** Unlike **Skippy** and LEAP agents, the Director agent receives  
 852 as input a simplified RGB image of the current state of the environment  
 853 (see Fig. 8). This is because we found that Director performed better  
 854 with its original architecture, which was designed for image-based  
 855 observations. We removed all textures to simplify the RGB observations.

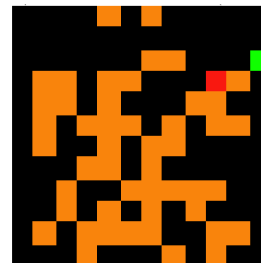


Figure 8: **An example for simplified observations for Director.**

856 **E.3.2 Failure**  
 857 **Modes: Bad Generalization, Sensitive to Short Trajectories**

858 **Training Performance.** We investigated why Director is unable to achieve good training performance (Fig. 3). As Director was designed to be trained solely on environments where it is able to collect long trajectories to train a good enough recurrent world model [23], we hypothesized that Director may perform better in domains where it is able to interact with the environment through longer trajectories by having better recurrent world models (*i.e.*, the agent does not immediately die as a result of interacting with specific objects in the environment). To test this, we experimented with variants of the used tasks, where the lava cells are replaced with wall cells, so the agent does not die upon trying to move towards them (we refer to this environment as the “walled” environment). The corresponding results on 50 training tasks are depicted in Fig. 9. As can be seen, the Director agent indeed performs better within the training tasks than in the environments with lava.

868 **Generalization Performance.** We also investigated why Director is unable to achieve good generalization (Fig. 3). As Director trains its policies solely from the imagined trajectories predicted by its learned world model, we believe that the low generalization performance is due to Director being unable to learn a good enough world model that generalizes to the evaluation tasks. The generalization performances in both the “walled” and regular environments, depicted in Fig. 9, indeed support this argument. Similar to what we did in the main paper, we also present experimental results for how the generalization performance changes with the number of training environments. Results in Fig. 10 show that the number of training environments has little effect on its poor generalization performance.

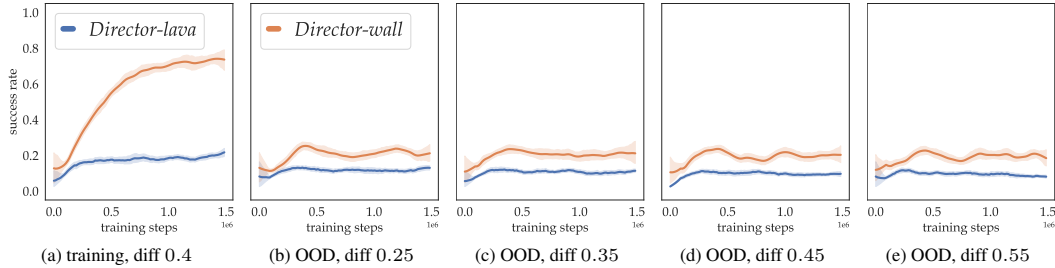


Figure 9: **Comparative Results of Director on Environments with Lavas and on those with Walls:** the results are obtained with 50 training tasks. The results for Director-lava (same as in the main paper) are obtained from 20 independent seed runs.

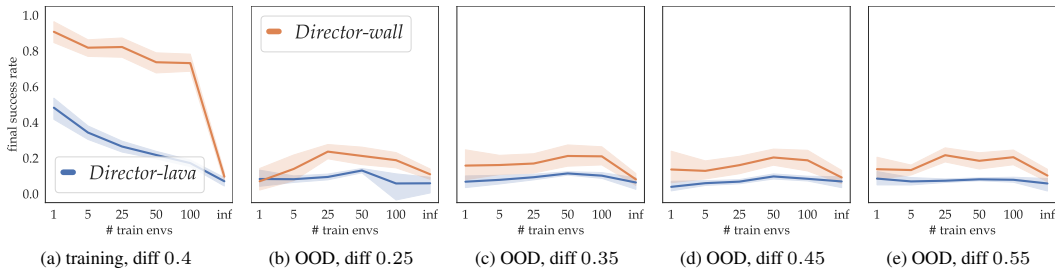


Figure 10: **Generalization Performance of Agents on Different Numbers of Training Tasks (while Director runs on the walled environments):** besides Director, each data point and corresponding error bar (95% confidence interval) are processed from the final performance from 20 independent seed runs. Director-wall’s results are obtained from 20 runs.

877 **F Experimental Results (Cont.)**

878 We present the experimental results that the main paper could not hold due to the page limit.

879 **F.1 Scalability of Generalization Performance**

880 Like [11], we investigate the scalability of the agents’ generalization abilities across different numbers  
 881 of training tasks. To this end, in Fig. 11, we present the results of the agents’ final evaluation  
 882 performance after training over different numbers of training tasks.

883 With more training tasks, **Skippers** and the baseline show consistent improvements in generalization  
 884 performance. While both LEAP and Director behave similarly as in the previous subsection, notably,  
 885 the **modelfree** baseline can reach similar performance as **Skippers**, but only when trained on a  
 886 different task in each episode, which is generally infeasible in the real world beyond simulation.

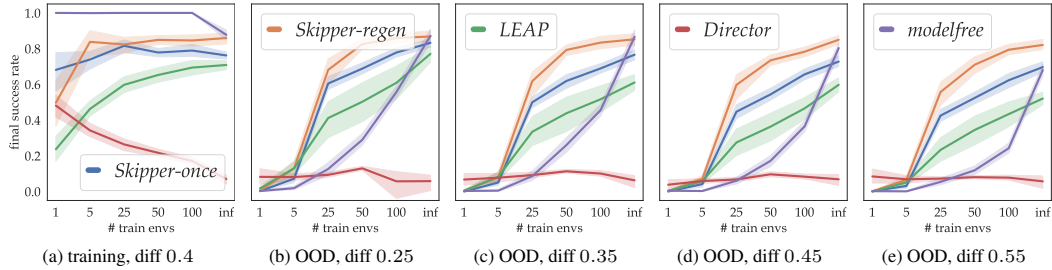


Figure 11: **Generalization Performance of Agents on Different Numbers of Training Tasks:** each data point and corresponding error bar (95% confidence interval) are based on the final performance from 20 independent seed runs. Training task performance is shown in (a) while OOD performance is shown in (b), (c), (d), (e). Notably, the **Skippers** agents as well as the adapted LEAP behave poorly during training when being trained on only one task, as the split of context and partial information cannot be achieved. Training on one task invalidates the purpose of the proposed generalization-focused checkpoint generator.

887 **F.2 Skipper-once Scalability**

888 We present the performance of **Skippers-once** on different numbers of training tasks in Fig. 12.

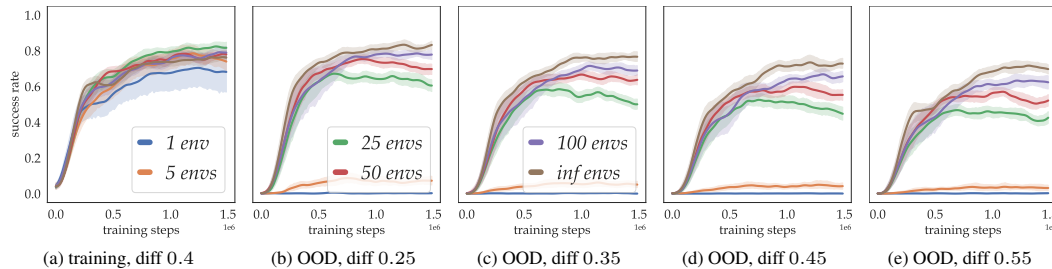


Figure 12: **Generalization Performance of Skipper-once on different numbers of training tasks:** each error bar (95% confidence interval) is obtained from 20 independent seed runs.

889 **F.3 Skipper-regen Scalability**

890 We present the performance of **Skippers-regen** on different numbers of training tasks in Fig. 13.

891 **F.4 modelfree Baseline Scalability**

892 We present the performance of the **modelfree** baseline on different numbers of training tasks in Fig.  
 893 14.

894 **F.5 LEAP Scalability**

895 We present the performance of the adapted LEAP baseline on different numbers of training tasks in  
 896 Fig. 15.



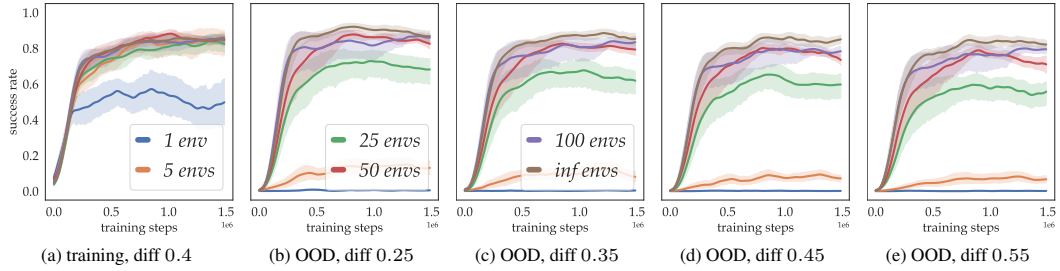


Figure 13: Performance of Skipper-regen on different numbers of training tasks: each error bar (95% confidence interval) is obtained from 20 independent seed runs.

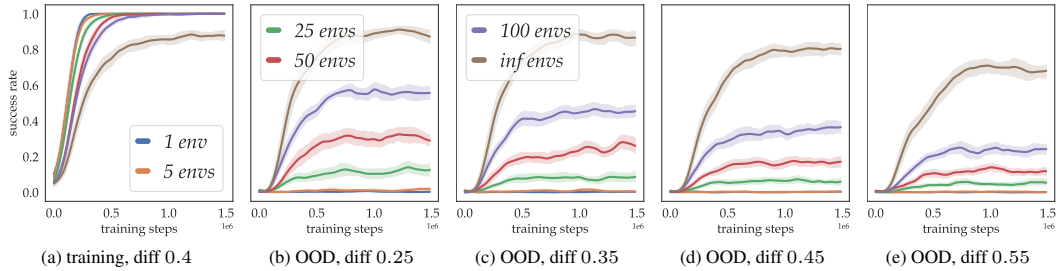


Figure 14: Generalization Performance of the modelfree baseline on different numbers of training tasks: each error bar (95% confidence interval) is obtained from 20 independent seed runs.

## 897 F.6 Director Scalability

898 We present the performance of the adapted Director baseline on different numbers of training tasks in  
 899 Fig. 16.

## 900 F.7 Generalization Performance on Different Numbers of Training Tasks

901 The performance of all agents on all training configurations, *i.e.* different numbers of training tasks,  
 902 are presented in Fig. 17, Fig. 18, Fig. 19, Fig. 20, Fig. 21 and Fig. 22.

### 903 F.7.1 Statistical Significance & Power Analyses

904 Besides visually observing generally non-overlapping confidence intervals, we present the pairwise  
 905 *t*-test results of Skipper-once and Skipper-regen against the compared methods. In addition, if the  
 906 advantage is significant, we perform power analyses to determine if the number of seed runs (20) was  
 907 enough to make the significance claim. These results are shown in Tab. 2 and Tab. 3, respectively.

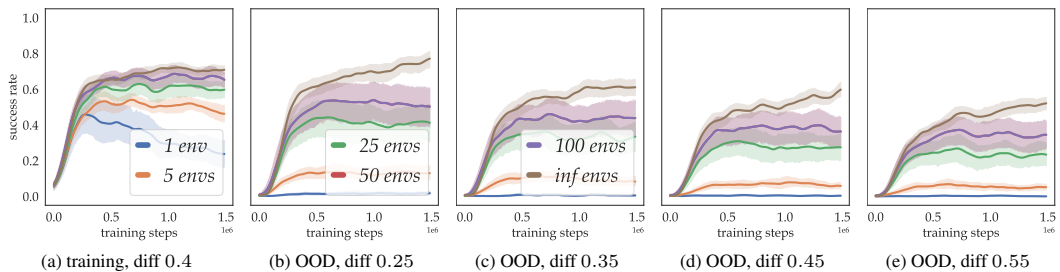


Figure 15: Generalization Performance of the LEAP baseline on different numbers of training tasks: each error bar (95% confidence interval) is obtained from 20 independent seed runs.

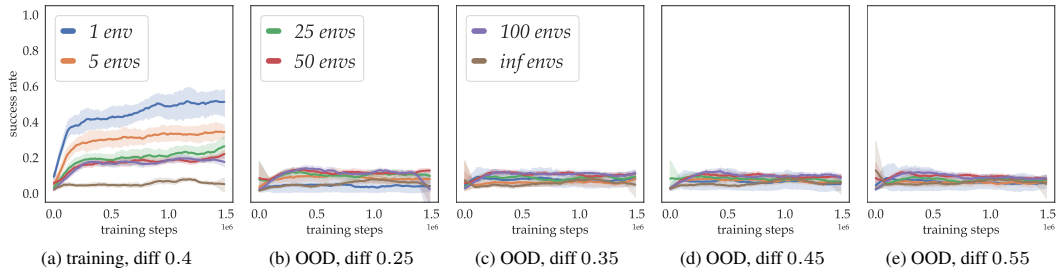


Figure 16: **Generalization Performance of the Director baseline on different numbers of training tasks:** each error bar (95% confidence interval) is obtained from 20 independent seed runs.

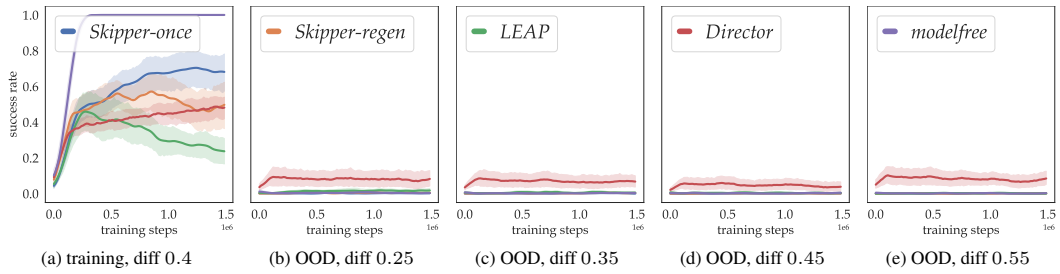


Figure 17: **Generalization Performance of the Agents when trained with 1 training task:** each error bar (95% confidence interval) is obtained from 20 independent seed runs.

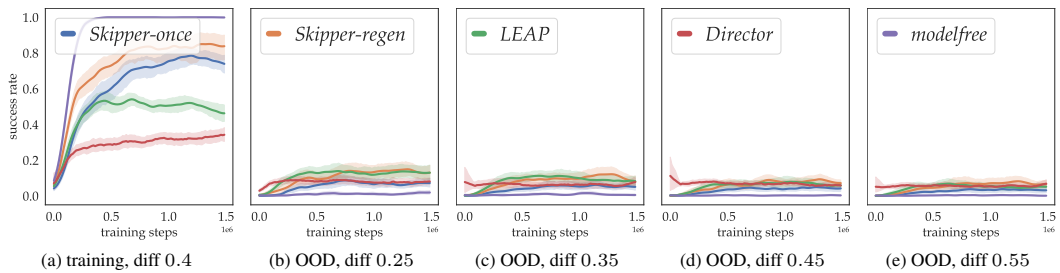


Figure 18: **Generalization Performance of the Agents when trained with 5 training tasks:** each error bar (95% confidence interval) is obtained from 20 independent seed runs.

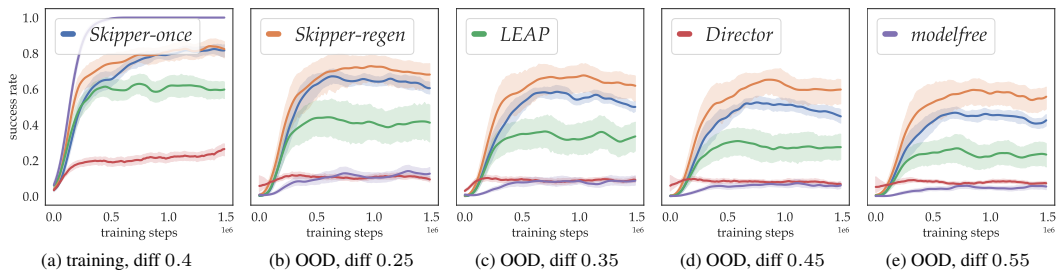


Figure 19: **Generalization Performance of the Agents when trained with 25 training tasks:** each error bar (95% confidence interval) is obtained from 20 independent seed runs.

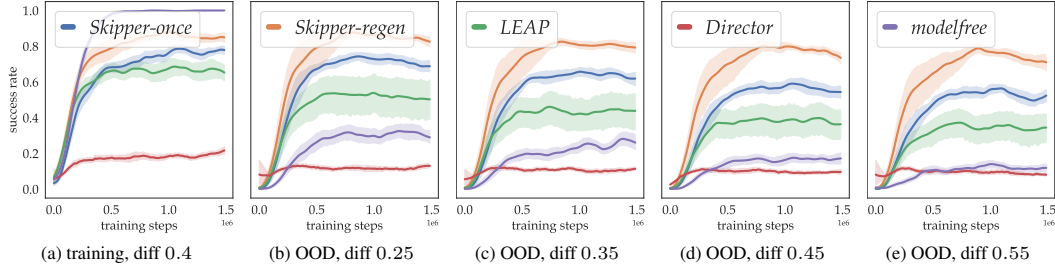


Figure 20: **Generalization Performance of the Agents when trained with 50 training tasks (same as in the main paper):** each error bar (95% confidence interval) is obtained from 20 independent seed runs.

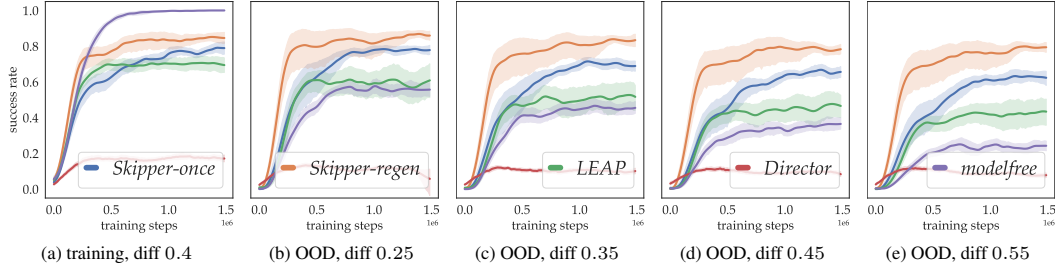


Figure 21: **Generalization Performance of the Agents when trained with 100 training tasks:** each error bar (95% confidence interval) is obtained from 20 independent seed runs.

908 As we can observe from the tables, generally there is significant evidence of generalization advantage  
 909 in **Skipper** variants compared to the other methods, especially when the number of training environ-  
 910 nments are between 25 to 100. Additionally, as expected, **Skipper-regen** displays more dominating  
 911 performance compared to that of **Skipper-once**.

## 912 G Ablation & Sensitivity

### 913 G.1 Validation of Effectiveness on Stochastic Environments

914 We present the performance of the agents in stochastic variants of the used environment. Specifically,  
 915 with probability 0.1, each action is changed into a random action. We present the 50-training  
 916 tasks performance evolution in Fig. 23. The results validate the compatibility of our agents with  
 917 stochasticity in environmental dynamics. Notably, the performance of the baseline deteriorated  
 918 to worse than even Director with the injected stochasticity. The compatibility of Hierarchical RL  
 919 frameworks to stochasticity has been investigated in [26].

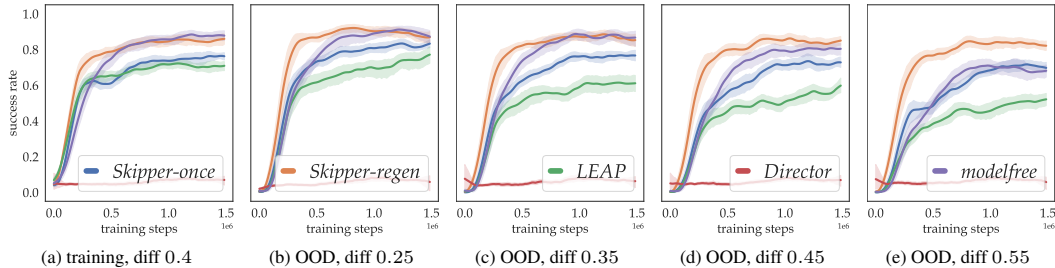


Figure 22: **Generalization Performance of the Agents when trained with  $\infty$  training tasks (a new task each training episode):** each error bar (95% confidence interval) is obtained from 20 independent seed runs.

Table 2: **Skipper-once** v.s. others: significance & power

	method \ task difficulty	0.25	0.35	0.45	0.55
1 train envs	leap	<b>22</b>	<b>NO</b>	<b>NO</b>	<b>NO</b>
	director	15	11	<b>22</b>	11
	baseline	<b>NO</b>	<b>38</b>	<b>36</b>	<b>NO</b>
5 train envs	leap	<b>28</b>	<b>NO</b>	<b>NO</b>	<b>NO</b>
	director	<b>NO</b>	<b>NO</b>	<b>NO</b>	<b>22</b>
	baseline	11	8	10	12
25 train envs	leap	15	13	11	7
	director	2	2	2	2
	baseline	2	2	2	2
50 train envs	leap	17	16	11	11
	director	2	2	2	2
	baseline	2	2	2	2
100 train envs	leap	15	10	7	9
	director	2	2	2	2
	baseline	2	2	2	2
inf train envs	leap	<b>32</b>	5	7	3
	director	2	2	2	2
	baseline	<b>NO</b>	<b>NO</b>	<b>NO</b>	<b>NO</b>

$t$  threshold: 0.05.

Effect size set to be the difference of the means of the compared pairs [12].

Cells are **bold** if results **NOT significant** or **insufficient seeds for statistical power**.

For significant cases, the minimum number of seeds for statistical power 0.2 is provided.

Table 3: **Skipper-regen** v.s. others: significance & power

	method \ task difficulty	0.25	0.35	0.45	0.55
1 train envs	leap	<b>32</b>	<b>NO</b>	<b>NO</b>	<b>NO</b>
	director	16	13	<b>23</b>	10
	baseline	<b>NO</b>	<b>NO</b>	<b>NO</b>	<b>NO</b>
5 train envs	leap	<b>NO</b>	<b>NO</b>	<b>NO</b>	<b>NO</b>
	director	<b>33</b>	<b>NO</b>	<b>NO</b>	<b>NO</b>
	baseline	6	8	4	5
25 train envs	leap	10	7	5	4
	director	2	2	2	2
	baseline	2	2	2	2
50 train envs	leap	6	4	3	3
	director	2	2	2	2
	baseline	2	2	2	2
100 train envs	leap	7	3	3	2
	director	2	2	2	2
	baseline	2	2	2	2
inf train envs	leap	15	3	2	2
	director	2	2	2	2
	baseline	<b>NO</b>	<b>NO</b>	<b>35</b>	5

$t$  threshold: 0.05.

Effect size set to be the difference of the means of the compared pairs [12].

Cells are **bold** if results **NOT significant** or **insufficient seeds for statistical power**.

For significant cases, the minimum number of seeds for statistical power 0.2 is provided.

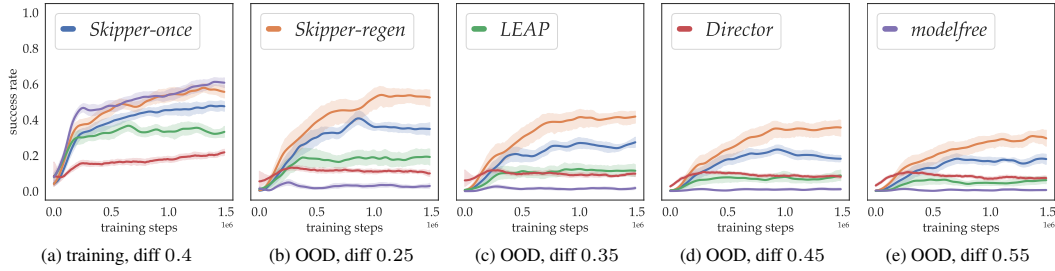


Figure 23: **Generalization Performance of agents in stochastic environments:**  $\epsilon$ -greedy style randomness is added to each primitive action with  $\epsilon = 0.1$ . Each agent is trained with 50 environments and each curve is processed from 20 independent seed runs.

920 **G.2 Ablation for Spatial Abstraction**

921 We present in Fig. 24 the ablation results on the spatial abstraction component with **Skipper-*once***  
 922 agent, trained with 50 tasks. The alternative component of the attention-based bottleneck, which is  
 923 without the spatial abstraction, is an MLP on a flattened full state. The results confirm significant  
 924 advantage in terms of generalization performance as well as sample efficiency in training, introduced  
 925 by spatial abstraction.

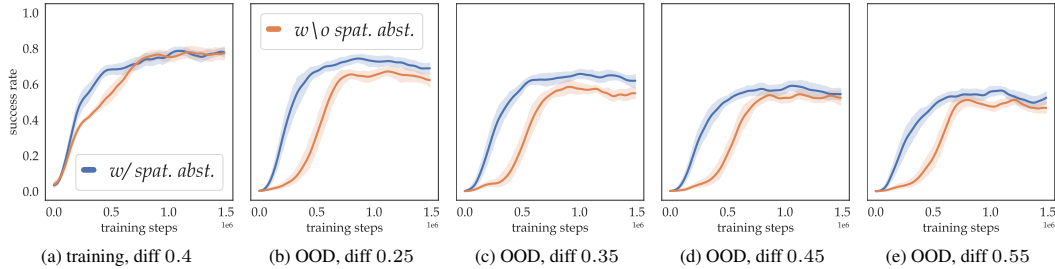


Figure 24: **Ablation for Spatial Abstraction on Skipper-*once* agent:** each agent is trained with 50 environments and each curve is processed from 20 independent seed runs.

926 **G.3 Accuracy of Proxy Problems & Checkpoint Policies**

927 We present in Fig. 25 the ablation results on the accuracy of proxy problems as well as the checkpoint  
 928 policies of the **Skipper-*once*** agents, trained with 50 tasks. The ground truths are computed via DP  
 929 on the optimal policies, which are also suggested by DP. Concurring with our theoretical analyses,  
 930 the results indicate that the performance of **Skipper** is determined (bottlenecked) by the accuracy of  
 931 the proxy problem estimation on the high-level and the optimality of the checkpoint policy on the  
 932 lower level. Specifically, the curves for the generalization performance across training tasks, as in (a)  
 933 of 25, indicate that the lower than expected performance is a composite outcome of errors in the two  
 934 levels. In the next part, we address a major misbehavior of inaccurate proxy problem estimation -  
 935 chasing delusional targets.

936 **G.4 Training Initialization: uniform v.s. same as evaluation**

937 We compare the agents' performance with and without uniform initial state distribution. The non-  
 938 uniform starting state distributions introduce additional difficulties in terms of exploration. In  
 939 Presented in Fig. 26, these results are obtained from training on 50 tasks. We conclude that given  
 940 similar computational budget, using non-uniform initialization only slows down the learning curves  
 941 without introducing significant changes to our conclusions, and thus we use the ones with uniform  
 942 initialization for presentation in the main paper.

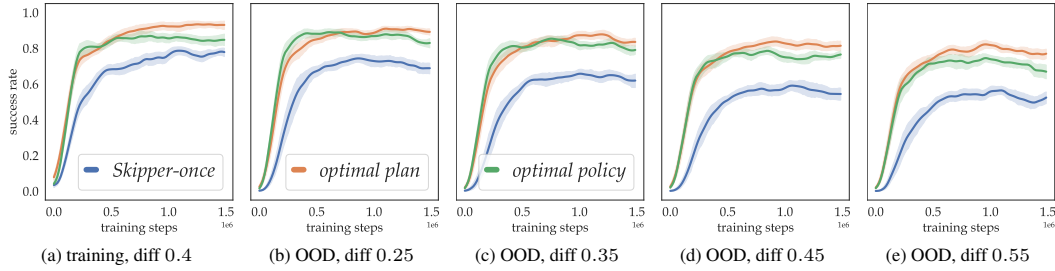


Figure 25: **Skipper-once Empirical Performance v.s. ground truths**: both the optimal policy and optimal plan variants are assisted by DP. The default deterministic setting induces the fact that combining optimal policy and optimal plan results in 1.0 success rate. The figures suggest that the learned agent is limited by errors both in the proxy problem estimation and the checkpoint policy  $\pi$ . Each agent is trained with 50 environments and each curve is processed from 20 independent seed runs.

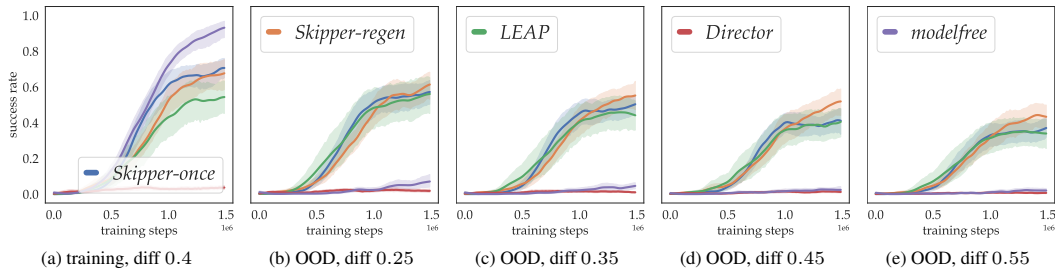


Figure 26: **Comparative Results on 50 training tasks without uniform initial state distribution**: each curve is processed from 20 independent seed runs.

## 943 G.5 Ablation: Vertex Pruning

944 As mentioned previously, each proxy problem in the experiments are reduced from 32 vertices to  
 945 12 with such techniques. We compare the performance curves of the used configuration against  
 946 a baseline that generates 12-vertex proxy problems without pruning. We present in Fig. 27 these  
 947 ablation results on the component of  $k$ -medoids checkpoint pruning. We observe that the pruning not  
 948 only increases the generalization but also the stability of performance.

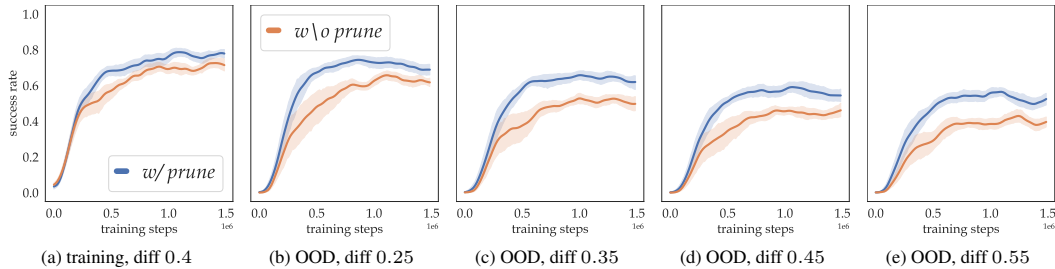


Figure 27: **Ablation Results on 50 training tasks for  $k$ -medoids pruning**: each curve is processed from 20 independent seed runs.

## 949 G.6 Sensitivity: Number of Vertices

950 We provide a sensitivity analysis to the number of checkpoints (number of vertices) in each proxy  
 951 problem. We present the results of **Skipper-once** on 50 training tasks with different numbers of  
 952 post-pruning checkpoints (all reduced from 32 by pruning), in Fig. 28. From the results, we can  
 953 see that as long as the number of checkpoints is above 6, **Skipper** exhibits good performance. We  
 954 therefore chose 12, the one with a rather small computation cost, as the default hyperparameter.

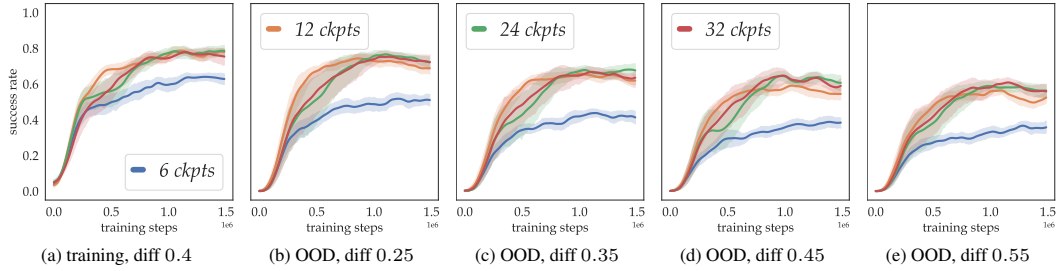


Figure 28: Sensitivity of Skipper-once on the number of checkpoints in each proxy problem: each agent is trained with 50 environments. All curves are processed from 20 independent seed runs.

## 955 G.7 Ablation: Planning over Proxy Problems

956 We provide additional results for the readers to intuitively understand the effectiveness of planning  
 957 over proxy problems. This is done by comparing the results of **Skipper-once** with a baseline **Skipper-**  
 958 **goal** that blindly selects the task goal as its target all the time. We present the results based on 50  
 959 training tasks in Fig. 29. Concurring with our vision on temporal abstraction, we can see that solving  
 960 more manageable sub-problems leads to faster convergence. The **Skipper**-goal variant catches up  
 961 later when the policy slowly improves to be capable of solving longer distance navigation.

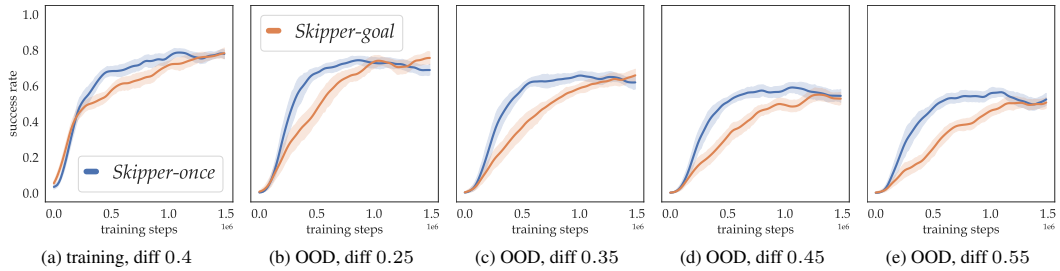


Figure 29: Effectiveness of Proxy Problem based Planning: each agent is trained with 50 environments and each curve is processed from 20 independent seed runs.