
World Models Increase Autonomy in Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Reinforcement learning (RL) is an appealing paradigm for training intelligent
2 agents, enabling policy acquisition from the agent’s own autonomously acquired
3 experience. However, the training process of RL is far from automatic, requir-
4 ing extensive human effort to reset the agent and environments. To tackle the
5 challenging reset-free setting, we first demonstrate the superiority of model-based
6 (MB) RL methods in such setting, showing that a straightforward application
7 of MBRL can outperform all the prior state-of-the-art methods while requiring
8 less supervision. We then identify limitations inherent to this direct extension
9 and propose a solution called model-based reset-free (MoReFree) agent, which
10 further enhances the performance. MoReFree adapts two key mechanisms, explo-
11 ration and policy learning, to handle reset-free tasks by prioritizing task-relevant
12 states. It exhibits superior data-efficiency across various reset-free tasks without
13 access to environmental reward or demonstrations while significantly outperform-
14 ing privileged baselines that require supervision. Our findings suggest model-
15 based methods hold significant promise for reducing human effort in RL. Website:
16 <https://sites.google.com/view/morefree>

17 1 Introduction

18 Reinforcement learning presents an attractive framework for training capable agents. At first glance,
19 RL training appears intuitive and autonomous - once a reward is defined, the agent learns from its
20 own automatically gathered experience. However, in practice, RL training often assumes the access
21 to environmental resets that can require significant human effort to setup, which poses a significant
22 barrier for real world applications of RL like robotics.

23 Most RL systems on real robots to date have employed various strategies to implement resets, all
24 requiring a considerable amount of effort [17, 32, 34, 21]. In [21], which trains a dexterous hand to
25 rotate balls, the practitioners had to (1) position a funnel underneath the hand to catch dropped balls,
26 and (2) deploy a separate robot arm to pick up the dropped balls for resets, and (3) script the reset
27 behavior. These illustrate that even for simple behaviors, proper implementation of reset mechanisms
28 can result in significant human effort and time.

29 Rather than depending on human-engineered reset mechanisms, the agent can operate within a reset-
30 free training scheme, learning to reset itself [4, 27, 25, 12] or train a policy capable of starting from
31 diverse starting states [35]. However, the absence of resets introduces unique exploration challenges.
32 Without periodic resets, the agent can squander significant time in task-irrelevant regions that require
33 careful movements to escape and may overexplore, never returning from indefinite exploration.
34 Recent unsupervised model-based RL (MBRL) approaches [20, 14] in the episodic setting have
35 shown sophisticated exploration, high data-efficiency and promising results in long-horizon tasks.
36 This prompts the question: *would MBRL agents excel in the reset-free RL setting?*

37 As an initial attempt, we first evaluate an unsupervised
 38 MBRL agent, out-of-the-box, in a reset-free Ant locomotion
 39 task. The ant is reset to the center of a rectangular
 40 arena, and is tasked with navigating to the upper right
 41 corner. The agent is reset only once at the start of training.
 42 The evaluation is episodic - the agent is reset at the start
 43 of each evaluation episode.

44 For the MBRL agent, we use PEG [14], which was devel-
 45 oped to solve hard exploration tasks in the episodic setting.
 46 As seen in Figure 1, PEG, out of the box, outperforms
 47 prior state-of-the-art, model-free agent, IBC [15], tailored
 48 for the reset-free setting.

49 In Figure 1, we plot state visitation heatmaps of the agents,
 50 where lighter colors correspond to more visitations. The
 51 oracle agent, with access to resets, explores the the “task-
 52 relevant” area between the initial and top right corner,
 53 which is ideal for training a policy that succeeds in episodic evaluation.
 54 IBC’s heatmap (bottom) shows that it fails to explore effectively, never encountering the goal states in the top right region.
 55 In contrast, PEG exhaustively explores the entire space, as seen through its uniform heatmap. This
 56 results in an overexploration problem - PEG may devote considerable time on finding irrelevant states
 57 rather than concentrating on the task-relevant region of the task. This leads us to ask: **how can MBRL**
 58 **agents acquire more task-relevant data in the reset-free setting to improve its performance?**

59 We propose **Model-based, Reset-Free (MoReFree)**, which improves two key mechanisms in model-
 60 based RL, exploration and policy optimization, to better handle reset-free training. Following the
 61 top row of Figure 2: to gather task-relevant data without resets, we define a training curriculum that
 62 alternates between temporally extended phases of task solving, resetting, and exploration. Next,
 63 as seen in the bottom row of Figure 2, we bias the policy training within the world model towards
 64 achieving task-relevant goals such as reaching initial states and evaluation states.

65 Our key contributions are as follows: **(1)** We demonstrate the viability of using model-based agents
 66 with strong exploration abilities for the reset-free setting as well as their inherent limitations. We
 67 address such limitations through the MoReFree framework which focuses exploration and policy opti-
 68 mization on task-relevant states. **(2)** We evaluate the out-of-the-box MBRL baseline and MoReFree
 69 against state-of-the-art reset-free methods in 8 challenging reset-free tasks ranging from manipulation
 70 to locomotion. Notably, both model-based approaches outperform prior state-of-the-art baselines
 71 in 7/8 tasks in final performance and data efficiency, all the while requiring less supervision (e.g.
 72 environmental reward or demonstrations). MoReFree outperforms the model-based baseline in the
 73 3 hardest tasks. **(3)** We perform in-depth analysis of the MoReFree and baselines behaviors, and
 74 show that MoReFree explores the state space thoroughly while retaining high visitation counts in the
 75 task-relevant regions. Our ablations show that the performance gains of MoReFree come from the
 76 proposed design choices and justify the approach.

77 2 Related Work

78 **Reset-free RL:** There is a growing interest in researching reinforcement learning methods that can
 79 effectively address the complexities of reset-free training. [28] proposes a reset-free RL benchmark
 80 (EARL) and finds that standard RL methods like SAC [9] fail catastrophically in EARL. Multiple
 81 approaches have been proposed to address reset-free training, which we now summarize. One
 82 approach is to add an additional reset policy, to bring the agent back to suitable states for learning
 83 [4, 16, 27, 25, 15]. LNT [4] and [16] train a reset policy to bring the agent back to initial state
 84 distribution, supervised by dense rewards and demonstrations respectively. MEDAL [25, 26], train a
 85 goal-conditioned reset policy and direct it to reset goal states from demonstrations. IBC [15] defines
 86 a curriculum for both task and reset policies without requiring demonstrations. VaPRL [27] trains a
 87 single goal-conditioned policy to reach high value states close to the initial states. Instead of guiding
 88 the agent back to familiar states, R3L [35] and [31] learn to reset the policy to diverse initial states,
 89 resulting in a policy that is more robust to variations in starting states. However, such methods are

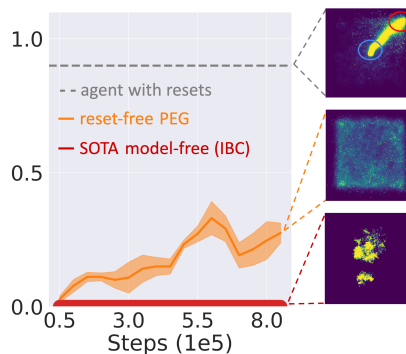


Figure 1: Performance and collected data of different agents on the reset-free Ant locomotion task.

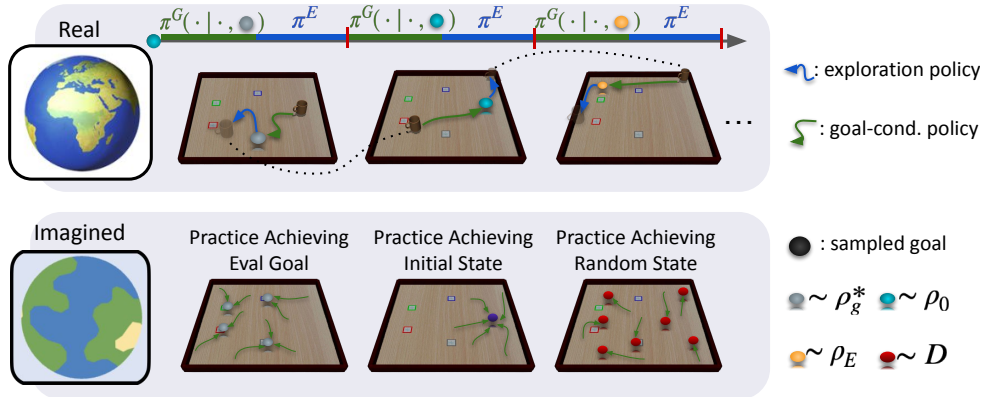


Figure 2: MoReFree is a model-based RL agent for solving reset-free tasks. **Top row:** MoReFree strikes a balance between exploring unseen states and practicing optimal behavior in task-relevant regions by directing the goal-conditioned policy to achieve evaluation states, initial state states (emulating a reset), and exploratory goals. **Bottom row:** MoReFree focuses the goal-conditioned policy training inside the world model on achieving evaluation states, initial states, and random replay buffer states to better prepare the policy for the aforementioned exploration scheme.

90 limited to tasks where exploration is unchallenging. The vast majority of reset-free approaches are
 91 model-free, with a few exceptions [19, 18]. Other works [8, 29] model the reset-free RL training
 92 process as a multi-task RL problem and require careful definition of the task distribution such that the
 93 tasks reset each other.

94 **Goal-conditioned Exploration:** A common theme running through the aforementioned work is
 95 the instantiation of a curriculum, often through commanding goal-conditioned policies, to keep the
 96 agent in task-relevant portions of the environment while exploring. Closely related is the subfield
 97 of goal-conditioned exploration in RL, where a goal-conditioned agent selects its own goals during
 98 training time to generate data. There is a large variety of approaches for goal selection, such as task
 99 progress [1, 30], intermediate difficulty [6], value disagreement [33], state novelty [23, 22], world
 100 model error [14, 24], and more. Many goal-conditioned exploration methods use the “Go-Explore”
 101 [3] strategy, which first selects a goal and runs the goal-conditioned policy (“Go”-phase), and then
 102 switches to an exploration policy for the latter half of the episode (“Explore”-phase). PEG [14],
 103 which MoReFree uses, extends Go-Explore to the model-based setting, and utilizes the world model
 104 to plan states with higher exploration value as goals. However, such methods are not designed for the
 105 reset-free RL setting, and may suffer from *over-exploration* of task-irrelevant states.

Table 1: A conceptual overview of reset-free methods. Existing methods are model-free, and most of them require other forms of supervision (environmental reward or demonstrations or both). In performance, MoReFree improves over reset-free PEG, which significantly outperforms privileged baselines IBC, MEDAL and R3L.

Approach	MEDAL	IBC	VaPRL	R3L	reset-free PEG	MoReFree
Model-based	\times	\times	\times	\times	\checkmark	\checkmark
Other Supervision	\checkmark	\checkmark	\checkmark	\times	\times	\times

106 We notice that the majority of all prior work are model-free and may suffer from poor sample efficiency
 107 and exploration issues. In contrast, our model-based approaches use world models to efficiently train
 108 policies and perform non-trivial goal-conditioned exploration with minimal supervision. See Table 1
 109 for a conceptual comparison between prior work and two model-based methods (MoReFree and
 110 reset-free PEG).

111 3 Preliminaries

112 3.1 Reset-free RL

113 We follow the definition of reset-free RL from EARL [28], and extend it to the goal-
 114 conditioned RL setting. Consider the goal-conditioned Markov decision process (MDP) $\mathcal{M} =$
 115 $(\mathcal{S}, \mathcal{G}, \mathcal{A}, p, r, \rho_0, \rho_{g^*}, \gamma)$. At each time step t in the state $s_t \in \mathcal{S}$, a goal-conditioned policy $\pi(\cdot|s_t, g)$
 116 under the goal command $g \in \mathcal{G}$ selects an action $a_t \in \mathcal{A}$ and transitions to the next state s_{t+1} with
 117 the probability $p(s_{t+1}|s_t, a_t)$, and gets a reward $r(s_t, a_t, g)$. ρ_0 is the initial state distribution, ρ_{g^*} is
 118 the evaluation goal distribution, and γ is the discount factor.

119 The learning algorithm \mathbb{A} is defined: $\{s_i, a_i, s_{i+1}\}_{i=0}^{t-1} \mapsto (a_t, \pi_t)$, which maps the transitions
 120 collected until the time step t to the action a_t the agent should take in the non-episodic training and
 121 the best guess π_t of the optimal policy π^* on the evaluation goal distribution (ρ_{g^*}). In reset-free
 122 training the agent will only be reset to the initial state $s_0 \sim \rho_0$ once. The evaluation of agents is still
 123 episodic. The agent always starts from $s_0 \sim \rho_0$, and is asked to achieve $g \sim \rho_{g^*}$. The evaluation
 124 objective for a policy π is:

$$J(\pi) = \mathbb{E}_{s_0 \sim \rho_0, g \sim \rho_{g^*}, a_j \sim \pi(\cdot|s_j, g), s_{j+1} \sim p(\cdot|s_j, a_j)} \left[\sum_{j=0}^T \gamma^j r(s_j, a_j, g) \right], \quad (1)$$

125 where T is the total time steps during the evaluation. The goal of algorithm \mathbb{A} during the reset-free
 126 training is to minimize the performance difference $\mathbb{D}(\mathbb{A})$ of the current policy π_t and the optimal
 127 policy π^* :

$$\mathbb{D}(\mathbb{A}) = \sum_{t=0}^{\infty} (J(\pi^*) - J(\pi_t)). \quad (2)$$

128 In summary, the algorithm \mathbb{A} should output an action a_t that the agent should take in the non-episodic
 129 data collection and a policy π_t that can maximize $J(\pi_t)$ at every time step t based on all previously
 130 collected data.

131 3.2 Model-based RL setup

132 Recent goal-conditioned MBRL approaches like LEXA [20] and PEG [14] train goal-conditioned
 133 policies purely using synthetic data generated by learned world models. Their robust exploration
 134 demonstrate significant success in solving long-horizon goal-conditioned tasks. In the reset-free
 135 setting, strong exploration is crucial, as the agent can no longer depend on episodic resets to bring it
 136 back to task-relevant areas if it gets stuck. Therefore, we select PEG as the backbone MBRL agent
 137 for its strong exploration abilities and sample efficiency.

138 PEG [14] is a model-based Go-Explore framework that extends LEXA [20], an unsupervised goal-
 139 conditioned variant of DreamerV2 [11]. The following components are parameterized by θ and
 140 learned:

$$\begin{aligned} \text{world model: } & \widehat{\mathcal{T}}_{\theta}(s_t|s_{t-1}, a_{t-1}) \\ \text{goal conditioned policy: } & \pi_{\theta}^G(a_t|s_t, g) & \text{goal conditioned value: } & V_{\theta}^G(s_t, g) \\ \text{exploration policy: } & \pi_{\theta}^E(a_t|s_t) & \text{exploration value: } & V_{\theta}^E(s_t) \end{aligned} \quad (3)$$

141 The world model is a recurrent state-space model (RSSM) which is trained
 142 to predict future states and is used as a learned simulator to train the
 143 policies and value functions. The goal-conditioned policy π^G is trained
 144 to reach random states sampled from the replay buffer. The exploration
 145 policy π^E is trained on an intrinsic motivation reward that rewards world
 146 model error, expressed through the variance of an ensemble [24]. Both
 147 policies are trained on simulated trajectory rollouts in the world model.

148 ► Self-supervised Goal-reaching Reward Function:

149 Rather than assuming access to the environmental reward, PEG learns
 150 its own reward function. PEG uses a dynamical distance function [13] as
 151 the reward function within world models, which predicts the number of

Algorithm 1 Go-Explore

```

1: Input:  $g, \pi_{\theta}^G, \pi_{\theta}^E$ 
2:  $\tau_g \leftarrow \{\}; \tau_e \leftarrow \{\}$ 
3: for  $t = 1$  to  $H_G$  do:
4:    $a_t \sim \pi^G(\cdot | s_t, g)$ 
5:    $s_{t+1} \sim \mathcal{T}(\cdot | s_t, a_t)$ 
6:    $\tau_g \leftarrow \tau_g \cup \{s_t\}$ 
7: for  $t = 1$  to  $H_E$  do:
8:    $a_t \sim \pi^E(\cdot | s_t)$ 
9:    $s_{t+1} \sim \mathcal{T}(\cdot | s_t, a_t)$ 
10:   $\tau_e \leftarrow \tau_e \cup \{s_t\}$ 
11: return  $\tau_g, \tau_e$ 

```

actions between a start and goal state. The distance function is trained on random state pairs from imaginary rollouts of π^G . π^G is then trained to minimize the dynamical distance between its states and commanded goal state in imagination. See [20] for more details.

► **Phased Exploration via Go-Explore:** For data-collection, PEG employs the Go-Explore strategy. In the “Go”-phase, a goal is sampled from some goal distribution ρ . The goal-conditioned policy, conditioned on the goal is run for some time horizon H_G , resulting in trajectory τ^G .

Then, in the “Explore”-phase, starting from the last state in the “Go”-phase, the exploration policy is run for H_E steps, resulting in τ^E . The interleaving of goal-conditioned behavior with exploratory behavior results in more directed exploration and informative data. This in turn improves accuracy of the world model, and the policies that train inside the world model. See

Algorithm 2 MBRL Backbone (PEG)

1: **Input:** $\pi_\theta^G, \pi_\theta^E$, world model \widehat{T}_θ , goal distribution ρ
2: **for** episode $i = 1$ to N **do**:
3: sample a goal $g \sim \rho$
4: $\tau_g, \tau_e \leftarrow \text{Go-Explore}(g, \pi^G, \pi^E)$
5: $\mathcal{D} \leftarrow \mathcal{D} \cup \tau^g \cup \tau^e$
6: update \widehat{T}_θ with \mathcal{D}
7: update π_θ^G and π_θ^E with \widehat{T}_θ in imagination

Algorithm 1 and Algorithm 2 for pseudocode. The choice of goal distribution ρ is important for Go-Explore. In easier tasks, the evaluation goal distribution ρ_{g^*} may be sufficient. But in longer-horizon tasks, evaluation goals may be too hard to achieve. Instead, intermediate goals from an exploratory goal distribution ρ_E can help the agent explore. We choose PEG, which generates goals by planning through the world model to maximize exploration value (see [14] for details).

4 Method

As motivated in Section 1 and Figure 1, the direct application of PEG to the reset-free setting shows promising performance but suffers from over-exploration of task-irrelevant states. To adapt model-based RL to the reset-free setting, we introduce MoReFree, a model-based approach that improves PEG to handle the lack of resets and overcome the over-exploration problem. MoReFree improves two key mechanisms of MBRL for reset-free training: exploration and policy training.

4.1 Back-and-Forth Go-Explore

First, we introduce MoReFree’s procedure for collecting new datapoints in the real environment. PEG [14] already has strong goal-conditioned exploration abilities, but was developed for solving episodic tasks. Without resets, PEG’s Go-Explore procedure can undesirably linger in unfamiliar but task-irrelevant portions of the state space. This generates large amounts of uninformative trajectories, which in turn degrades world model learning and policy optimization.

MoReFree overcomes this by periodically directing the agent to return to the states relevant to the task (i.e. initial and evaluation goals). We call this exploration procedure “Back-and-Forth Go-Explore”, where we sample pairs of initial and evaluation goals and ask the agent to cycle back and forth between the goal pairs, periodically interspersed with exploration phases (see Figure 2 top row).

Now, we define the “Back-and-Forth Go-Explore” strategy as seen in Algorithm 3. First, we decide whether to perform initial/evaluation state directed exploration. With probability α , we sample goals (g^*, g_0) from ρ_{g^*}, ρ_0 respectively. Then, we execute the Go-Explore routine for each goal. We name Go-Explore trajectories conditioned on initial state goals as “Back” trajectories, and Go-Explore trajectories conditioned on evaluation goals as “Forward” trajectories. With probability $1 - \alpha$, we execute exploratory Go-Explore behavior by sampling exploratory goals from PEG.

Algorithm 3 Back-and-Forth Go-Explore

1: **Input:** $\pi_\theta^G, \pi_\theta^E$, world model \widehat{T}_θ , $\rho_{g^*}, \rho_0, \rho_E$
2: Generate a random number r in $[0, 1]$
3: **if** $r < \alpha$ **then**
4: $g^*, g_0 \sim \rho_{g^*}, \rho_0$
5: $\tau_{g^*}, \tau_e^1 \leftarrow \text{Go-Explore}(g^*, \pi^G, \pi^E)$
6: $\tau_{g_0}, \tau_e^2 \leftarrow \text{Go-Explore}(g_0, \pi^G, \pi^E)$
7: **else**
8: $g \sim \rho_E$
9: $\tau_g, \tau_e^1 \leftarrow \text{Go-Explore}(g, \pi^G, \pi^E)$
10: **end if**

By following this exploration strategy, the agent modulates between various Go-Explore strategies, alternating between solving the task by pursuing evaluation goals, resetting the task by pursuing initial states, and exploring unfamiliar regions via exploratory goals.

204 4.2 Learning to Achieve Relevant Goals in Imagination

205 Next, we describe how MoReFree trains the goal-conditioned policy in the world model. To train
206 π^G , MoReFree samples various types of goals and executes $\pi^G(\cdot | \cdot, g)$ inside the world model to
207 generate “imaginary” trajectories. The trajectory data is scored using the learned dynamical distance
208 reward mentioned in Section 3.2, and the policy is updated to maximize the expected return. This
209 procedure is called imagination [10], and allows the policy to be trained on vast amounts of synthetic
210 trajectories to improve sample efficiency.

211 First, we choose to sample evaluation goals from ρ_{g^*} since the policy will be evaluated on its
212 evaluation goal-reaching performance. Next, recall that Back-and-Forth Go-Explore procedure
213 also samples initial states from ρ_0 as goals for the Go-phase to emulate resetting behavior. Since
214 we would like π^G to succeed in such cases so that the task is reset, we will also sample from ρ_0 .
215 Finally, we sample random states from the replay buffer to increase π^G ’s ability to reach arbitrary
216 states. The sampling probability for each goal type is set to $\alpha/2, \alpha/2, 1 - \alpha$ respectively. In other
217 words, MoReFree biases the goal-conditioned policy optimization procedure to focus on achieving
218 task-relevant goals (i.e. evaluation and initial states), as they are used during evaluation and goal-
219 conditioned exploration to condition the goal-reaching policy (see Figure 2 bottom row).

220 4.3 Implementation Details

221 Our work builds on the top of PEG [14], and we use its default hyperparameters for world model,
222 policies, value functions and temporal reward function. We set the length of each phase for Go-
223 Explore (H_G, H_E) to half the evaluation episode length for each task. We set the default value of
224 $\alpha = 0.2$ for all tasks (never tuned). See Appendix C.3 for more details and the supplemental for
225 MoReFree code.

226 5 Experiments

227 We evaluate two MBRL methods (PEG [14] and our extension MoReFree) and four competitive
228 reset-free baselines on eight reset-free tasks. We aim to address the following questions: 1) Do
229 MBRL approaches work well in reset-free tasks in terms of sample efficiency and performance? 2)
230 What limitations arise from running MBRL in the reset-free setting, and does our proposed solution
231 MoReFree address them? 3) What sorts of behavior do MoReFree and baselines exhibit in such tasks,
232 and are our design choices for MoReFree justified?

233 **Baselines:** All baselines except for R3L are implemented using official codebases, see Appendix C.2
234 for details.

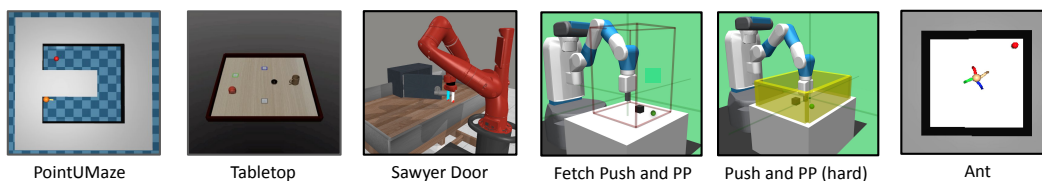


Figure 3: We evaluate MoReFree on eight reset-free tasks ranging from navigation to manipulation. PP is short for Pick&Place.

- 235 • **reset-free PEG** is a straightforward extension of PEG [14] to the reset-free setting.
- 236 • **MEDAL** [25] requires demonstrations and trains two policies, one for returning to demon-
237 stration states and another that achieves task goals.
- 238 • **IBC** [15] is a competitive baseline that outperforms prior reset-free work (e.g. MEDAL,
239 VaPRL) by defining a bidirectional curriculum for the goal-conditioned forward and back-
240 wards (i.e. reset) policies trained using the environmental reward.
- 241 • **R3L** [35] trains two policies, one for achieving task goals and another that perturbs the
242 agent to novel states. Notably, it is the only baseline that operates without any additional
243 assumptions (i.e. environmental rewards, demonstrations, and resets).

244

- **Oracle** is SAC [9] trained under the episodic setting on the environmental reward.

245

Note that most baselines enjoy some advantage over two MBRL methods: MEDAL, IBC and Oracle use ground truth environmental reward, while MEDAL also uses demonstrations and Oracle uses resets. See Table 1 for a conceptual comparison between MoReFree and prior work.

246

247

Environments: We evaluate MoReFree and baselines on eight tasks (see Figure 3). We select five tasks from IBC’s evaluation suite of six tasks; (Fetch Reach is omitted because it is trivially solvable). Next, we increased the complexity of the two hardest tasks from IBC, Fetch Push and Fetch Pick&Place, by extending the size of the workspace, replacing artificial workspace limits (which cause unrealistic jittering behavior near the limits, see the website for videos) with real walls, and evaluating on harder goal states (i.e. Pick&Place goals only in the air rather than including ones on the ground). In addition, we contributed a difficult locomotion task, Ant, which is adapted from the PEG codebase [14]. All methods are run with 5 seeds, and the mean performance and standard error are reported. During the evaluation, the performance on tasks with randomly sampled goals from ρ_{g^*} is measured by averaging over 10 episodes. See Appendix C for more experimental details.

257

258

5.1 Results

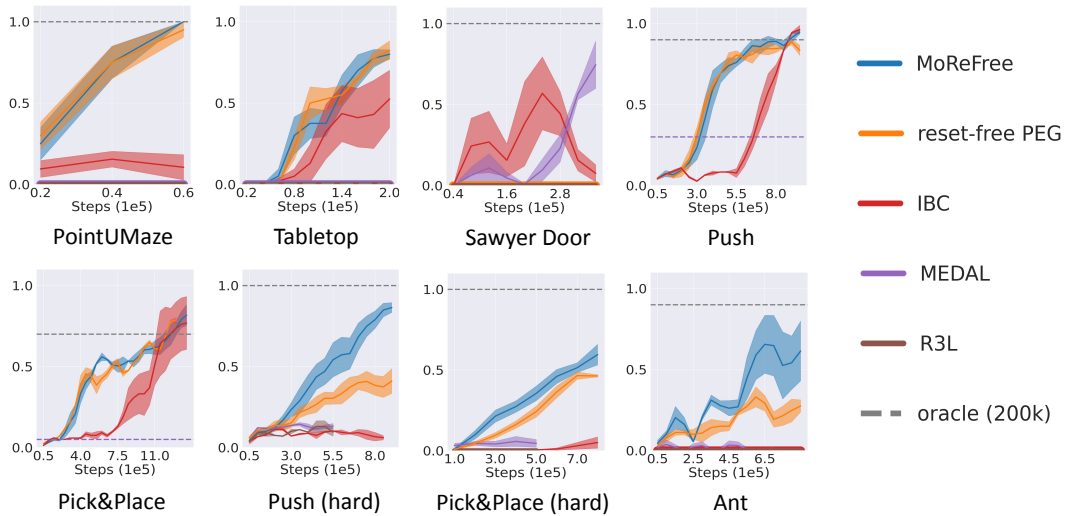


Figure 4: MBRL methods (MoReFree and reset-free PEG) significantly outperform baselines in 7/8 tasks. In 4 tasks, only MBRL methods are able to learn meaningful behavior, showcasing MBRL’s sample efficiency. MoReFree outperforms PEG in the 3 most difficult tasks.

259

As shown in Fig 4, two model-based methods (MoReFree and reset-free PEG), without demonstrations or access to environmental reward, outperform other baselines with privileged access to supervision in both final performance and sample efficiency in 7/8 tasks. We observe that the two MBRL methods learn good behaviors: the pointmass agent hugs the wall of the UMaze to minimize travel time and the Fetch robot deftly pushes and picks up the block into multiple target locations. MoReFree is always competitive with or outperforms PEG, with large gains in the 3 hardest tasks: Push (hard) by 45%, Pick&Place (hard) by 13% and Ant (hard) by 36%. We observe that MoReFree learns non-trivial reset behaviors such as picking and pushing blocks back into the center of the table for the hard variants of the Fetch manipulation tasks. See the website for videos of MoReFree and baselines.

267

268

In many tasks, the baselines fail to learn at all. We believe this is due the low sample budget, which may be too low for the baselines to fully explore the environment and learn the proper resetting behaviors necessary to train the actual task policy. In Appendix G, we increased the training budget by $3\times$ for the IBC baseline and it still fails, underscoring the difficulty of the tasks and the sample-efficiency gains of MoReFree and MBRL. On the other hand, we noticed that one environment, Sawyer Door, seemed particularly hard for MBRL agents to solve. We hypothesize that the dynamics of the task are hard to model, resulting in performance degradation for model-based approaches (see Appendix F for more analysis).

275

276 **5.2 Analysis**

277 To explain the performance differences between MoReFree and baselines, we closely analyze the
 278 exploration behaviors.

279 **MoReFree focuses on task-relevant states.** In Figure 5 we visualize the state visitation heatmaps of
 280 methods in various environments, and also compute the percentage of “task-relevant” states (initial
 281 and goal regions, highlighted with white borders). We highlight two trends. First, the heatmaps show
 282 that MoReFree and PEG explore thoroughly while baselines have more myopic exploration patterns,
 283 as seen in the Ant heatmaps at the top.

284 Next, performance differences between PEG and MoReFree are intuitively explained by the amount
 285 of task-relevant data collected by each agent. In easier environments like Push or Pick&Place
 286 where both PEG and MoReFree encounter similar amounts of task-relevant states, the performance
 287 is roughly similar between PEG and MoReFree. But in harder environments (Ant, Push (hard),
 288 Pick&Place (hard)) with larger state spaces and more complicated resetting dynamics, MoReFree
 289 collects $1.3 - 5\times$ more task-relevant data and has large performance gains over PEG. By experiencing
 290 more task-relevant states and training policies on them in imagination, MoReFree policies are
 291 more suited towards succeeding at the episodic evaluation criteria. See Appendix D for additional
 292 visualizations.

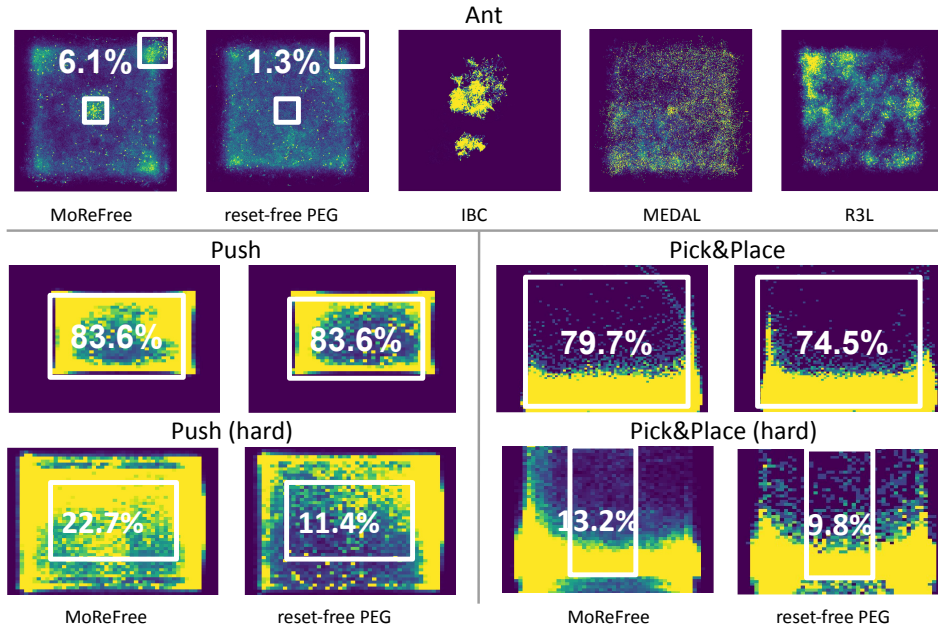


Figure 5: State visitation heatmaps of different agents. White areas are task-relevant states (including initial and goal state distributions) and we overlay the percentages of task-relevant states. MBRL methods explore more and in harder environments, MoReFree experiences more task-relevant states.

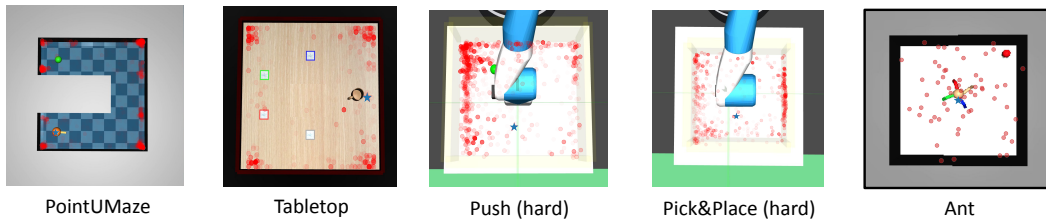


Figure 6: We visualize the start position (red dots) of successful “Back” trajectories of MoReFree’s Back-and-Forth Go-Explore, where π^G is directed to reset the environment. The color intensity of the dots correspond to state density over the last 100 steps.

293 **MoReFree effectively resets.** Next, we investigate the qualitative behavior of MoReFree’s Back-
 294 and-Forth Go-Explore. To see if “Back” trajectories help free the agent from the sink states, we
 295 analyze the replay buffer of MoReFree for the environments, and plot the starting locations of the
 296 agent / object up to 100 timesteps before a successful “Back” trajectory is executed in Figure 6. The
 297 color intensity of the dots correspond to state density over the last 100 steps (i.e. dark red means
 298 the agent / object has rested there for a while). We observe that the starting locations (red dots) of
 299 the agent / object are in corners or next to walls in all environments. This suggests that these areas
 300 act as sink states, where the agent / object would remain for long and waste time. We observe that
 301 MoReFree learns reset behaviors like picking the block out of corners and walls in Fetch Push and
 302 Fetch Pick&Place. See detailed videos of the reset behavior on the website.

303 5.3 Ablations

304 To justify our design choices, we ablate the two mecha-
 305 nisms of MoReFree, the back-and-forth exploration and
 306 goal-conditioned policy training, and plot the results in
 307 Figure 7 First, removing all mechanisms (**MF w/o Ex-
 308 plore & Imag.**) reduces to PEG, and we can see a large
 309 gap in performance. Next, **MF with Only Task Goals**
 310 sets $\alpha = 1$, which causes an extreme bias towards task-
 311 relevant states in the exploration and policy training. This
 312 also degrades performance, due to the need for strong ex-
 313 ploration in the reset-free setting. Examinations of more
 314 values for α can be found in Appendix C.3.

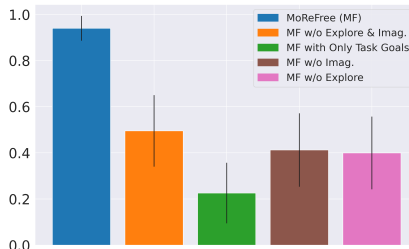


Figure 7: Ablations on 5 variants of MoReFree over 5 environments with normalized final performance.

315 Finally, we isolate individual components of MoReFree.
 316 First, we disable Back-and-Forth Go-Explore by disal-
 317 lowing the sampling of initial or evaluation goals during
 318 Go-Explore. Only exploratory goals are used in Go-Explore for this ablation (named **MF w/o
 319 BF-GE**). Next, in **MF w/o Imag.**, we turn off the initial / evaluation goal sampling in imagination, so
 320 only random replay buffer goals are used to train π^G . We see that both variants perform poorly. This
 321 is somewhat intuitive, as the two components rely on each other. Having both forms a synergistic
 322 cycle where 1) the goal-conditioned policy’s optimization is more focused towards reaching initial
 323 / goal states, and 2) the exploration is biased towards reaching initial / goal states by using the
 324 goal-conditioned policy we just optimized in step 1. If we remove one without the other, then the
 325 cycle breaks down. In **MF w/o Imag.**, Back-and-Forth Go-Explore will suffer since π^G trained on
 326 random goals cannot reliably reach initial / evaluation goals. In **MF w/o BF-GE**, the exploration
 327 strategy will not seek initial / evaluation states, resulting in an inaccurate world model and degraded
 328 policy optimization. In summary, the ablations show that MoReFree’s design is sound and is the
 329 major factor behind its success in the reset-free setting. See Appendix E for details.

330 6 Conclusion and Future Work

331 As a step towards reset-free training, we adapt model-based methods to the reset-free setting and
 332 demonstrate their superior performance. Specifically, we show that the out-of-the-box, unsupervised
 333 MBRL method substantially outperforms the state-of-the-art model-free baselines tailored for the
 334 reset-free setting while being more autonomous (requires less supervision like environmental reward
 335 or demonstrations). We then identify a limitation of unsupervised MBRL in the reset-free setting
 336 (over-exploration on task-irrelevant states), and propose MoReFree to address such limitations by
 337 focusing model-based exploration and goal-conditioned policy training on task-relevant states. We
 338 conduct a through experimental study of MoReFree and baselines over 8 tasks, and show considerable
 339 performance gains over the MBRL baseline and prior state-of-the-art reset-free methods. Despite
 340 its overall success, MoReFree is not without limitations. Being a model-based approach, it inherits
 341 all associated disadvantages. For example, we believe Sawyer Door is a task where learning the
 342 dynamics is harder than learning the policy (see Appendix F), disadvantaging MBRL approaches.
 343 Next, MoReFree uses a fixed percentage of task-relevant goals for exploration and imagination,
 344 whereas future work could consider an adaptive curriculum. Finally, scaling MoReFree to high-
 345 dimensional observations would be a natural extension. We hope MoReFree inspires future efforts in
 346 increasing autonomy in RL.

References

- 347
- 348 [1] A. Baranes and P.-Y. Oudeyer. Active learning of inverse models with intrinsically motivated
349 goal exploration in robots. *Robotics and Autonomous Systems*, 61(1):49–73, 2013.
- 350 [2] Y. Burda, H. Edwards, A. Storkey, and O. Klimov. Exploration by random network distillation.
351 *arXiv preprint arXiv:1810.12894*, 2018.
- 352 [3] A. Ecoffet, J. Huizinga, J. Lehman, K. O. Stanley, and J. Clune. First return, then explore.
353 *Nature*, 590(7847):580–586, 2021.
- 354 [4] B. Eysenbach, S. Gu, J. Ibarz, and S. Levine. Leave no trace: Learning to reset for safe and
355 autonomous reinforcement learning. *arXiv preprint arXiv:1711.06782*, 2017.
- 356 [5] B. Eysenbach, R. Salakhutdinov, and S. Levine. C-learning: Learning to achieve goals via
357 recursive classification. In *International Conference on Learning Representations*, 2021.
- 358 [6] C. Florensa, D. Held, X. Geng, and P. Abbeel. Automatic goal generation for reinforcement
359 learning agents. In *International conference on machine learning*, pages 1515–1528. PMLR,
360 2018.
- 361 [7] J. Fu, A. Singh, D. Ghosh, L. Yang, and S. Levine. Variational inverse control with events: A
362 general framework for data-driven reward definition. *Advances in neural information processing*
363 *systems*, 31, 2018.
- 364 [8] A. Gupta, J. Yu, T. Z. Zhao, V. Kumar, A. Rovinsky, K. Xu, T. Devlin, and S. Levine. Reset-free
365 reinforcement learning via multi-task learning: Learning dexterous manipulation behaviors with-
366 out human intervention. In *2021 IEEE International Conference on Robotics and Automation*
367 *(ICRA)*, pages 6664–6671. IEEE, 2021.
- 368 [9] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta,
369 P. Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*,
370 2018.
- 371 [10] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent
372 imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- 373 [11] D. Hafner, T. Lillicrap, M. Norouzi, and J. Ba. Mastering atari with discrete world models.
374 *arXiv preprint arXiv:2010.02193*, 2020.
- 375 [12] S. Haldar, J. Pari, A. Rai, and L. Pinto. Teach a robot to fish: Versatile imitation from one
376 minute of demonstrations. *arXiv preprint arXiv:2303.01497*, 2023.
- 377 [13] K. Hartikainen, X. Geng, T. Haarnoja, and S. Levine. Dynamical distance learning for semi-
378 supervised and unsupervised skill discovery. *arXiv preprint arXiv:1907.08225*, 2019.
- 379 [14] E. S. Hu, R. Chang, O. Rybkin, and D. Jayaraman. Planning goals for exploration. In *The*
380 *Eleventh International Conference on Learning Representations*, 2023.
- 381 [15] J. Kim, D. Cho, and H. J. Kim. Demonstration-free autonomous reinforcement learning via
382 implicit and bidirectional curriculum. In *International Conference on Machine Learning*. PMLR,
383 2023.
- 384 [16] J. Kim, J. hyeon Park, D. Cho, and H. J. Kim. Automating reinforcement learning with
385 example-based resets. *IEEE Robotics and Automation Letters*, 7(3):6606–6613, 2022.
- 386 [17] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies.
387 *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- 388 [18] K. Lu, A. Grover, P. Abbeel, and I. Mordatch. Reset-free lifelong learning with skill-space
389 planning. *arXiv preprint arXiv:2012.03548*, 2020.
- 390 [19] K. Lu, I. Mordatch, and P. Abbeel. Adaptive online planning for continual lifelong learning,
391 2020.

- 392 [20] R. Mendonca, O. Rybkin, K. Daniilidis, D. Hafner, and D. Pathak. Discovering and achieving
393 goals via world models, 2021.
- 394 [21] A. Nagabandi, K. Konolige, S. Levine, and V. Kumar. Deep dynamics models for learning
395 dexterous manipulation. In *Conference on Robot Learning*, pages 1101–1112. PMLR, 2020.
- 396 [22] S. Pitis, H. Chan, S. Zhao, B. Stadie, and J. Ba. Maximum entropy gain exploration for long
397 horizon multi-goal reinforcement learning. In *International Conference on Machine Learning*,
398 pages 7750–7761. PMLR, 2020.
- 399 [23] V. H. Pong, M. Dalal, S. Lin, A. Nair, S. Bahl, and S. Levine. Skew-fit: State-covering
400 self-supervised reinforcement learning. *arXiv preprint arXiv:1903.03698*, 2019.
- 401 [24] R. Sekar, O. Rybkin, K. Daniilidis, P. Abbeel, D. Hafner, and D. Pathak. Planning to explore
402 via self-supervised world models. In *ICML*, 2020.
- 403 [25] A. Sharma, R. Ahmad, and C. Finn. A state-distribution matching approach to non-episodic
404 reinforcement learning. *arXiv preprint arXiv:2205.05212*, 2022.
- 405 [26] A. Sharma, A. M. Ahmed, R. Ahmad, and C. Finn. Self-improving robots: End-to-end
406 autonomous visuomotor reinforcement learning. *arXiv preprint arXiv:2303.01488*, 2023.
- 407 [27] A. Sharma, A. Gupta, S. Levine, K. Hausman, and C. Finn. Autonomous reinforcement learning
408 via subgoal curricula. *Advances in Neural Information Processing Systems*, 34:18474–18486,
409 2021.
- 410 [28] A. Sharma, K. Xu, N. Sardana, A. Gupta, K. Hausman, S. Levine, and C. Finn. Autonomous
411 reinforcement learning: Formalism and benchmarking. *arXiv preprint arXiv:2112.09605*, 2021.
- 412 [29] L. Smith, N. Dhawan, M. Zhang, P. Abbeel, and S. Levine. Avid: Learning multi-stage tasks
413 via pixel-level translation of human videos. *arXiv preprint arXiv:1912.04443*, 2019.
- 414 [30] V. Veeriah, J. Oh, and S. Singh. Many-goals reinforcement learning. *ArXiv*, abs/1806.09605,
415 2018.
- 416 [31] K. Xu, S. Verma, C. Finn, and S. Levine. Continual learning of control primitives: Skill
417 discovery via reset-games. *Advances in Neural Information Processing Systems*, 33:4999–5010,
418 2020.
- 419 [32] A. Yahya, A. Li, M. Kalakrishnan, Y. Chebotar, and S. Levine. Collective robot reinforcement
420 learning with distributed asynchronous guided policy search. In *2017 IEEE/RSJ International
421 Conference on Intelligent Robots and Systems (IROS)*, pages 79–86. IEEE, 2017.
- 422 [33] Y. Zhang, P. Abbeel, and L. Pinto. Automatic curriculum learning through value disagreement.
423 *Advances in Neural Information Processing Systems*, 33:7648–7659, 2020.
- 424 [34] H. Zhu, A. Gupta, A. Rajeswaran, S. Levine, and V. Kumar. Dexterous manipulation with deep
425 reinforcement learning: Efficient, general, and low-cost. In *2019 International Conference on
426 Robotics and Automation (ICRA)*, pages 3651–3657. IEEE, 2019.
- 427 [35] H. Zhu, J. Yu, A. Gupta, D. Shah, K. Hartikainen, A. Singh, V. Kumar, and S. Levine. The
428 ingredients of real-world robotic reinforcement learning. *arXiv preprint arXiv:2004.12570*,
429 2020.

430 A Broader Impacts

431 As we increase the autonomy of RL agents, the possibility of them acting in unexpected ways to
432 maximize reward increases. The unsupervised exploration coupled alongside the learned reward
433 functions further add to the unpredictability; neither mechanisms are very interpretable. As such, we
434 expect research into value alignment, interpretability, and safety to be paramount as autonomy in RL
435 improves.

436 B Extended Related Work

437 **Learned Reward Functions:** Instead of requiring the environment to provide a reward function,
438 the agent can learn its own reward function from onboard sensors and data. Given human specified
439 example states, e.g. a goal image, VICE and C-Learning train reward classifiers over examples [7, 5]
440 and agent data. The learned dynamical distance function [13] learns to predict the number of actions
441 between pairs of states. The dynamical distance function is used by unsupervised MBRL approaches
442 like LEXA and PEG [20, 14] to train the goal-conditioned policy.

443 C Experimental Details

444 C.1 Environments

445 **PointUMaze:** The state space is 7D and the action space is 2D. The initial state is $(0, 0)$, which
446 located in the bottom-left corner, and noise sampled from $\mathcal{U}(-0.1, 0.1)$ is added when reset. The goal
447 during the evaluation is always located in at the top-left corner of the U-shape maze. The maximum
448 steps during the evaluation is 100. Hard reset will happen after every $2e5$ steps. In the whole training
449 process we performed, it only reset once at the beginning of the training. Taken from the IBC [15]
450 paper.

451 **Tabletop:** The state space is 6D, and the action space is 3D. During the evaluation, four goal locations
452 are sampled in turn, the initial state of the agent is always fixed and located in the center of the table.
453 The maximum steps during the evaluation is 200. Hard reset will happens after every $2e5$ steps. In
454 the whole training process we performed, it only reset once at the beginning of the training. Taken
455 from the EARL [28] benchmark and also used in the IBC paper.

456 **Sawyer Door:** The state space is 7D and the action space is 4D. The position of door is initialized to
457 open state (60 degree with noise sampled from $(0, 18)$ degree) and the goal is always to close the
458 door (0 degree). The arm is initialized to a fixed location. Maximum number of steps is 300 for the
459 evaluation. Hard reset will happen after every $2e5$ steps. In the whole training process we performed,
460 it resets twice. Taken from the EARL [28] benchmark and also used in the IBC paper.

461 **Fetch Push and Pick&Place:** The state space is 25D and action space is 4D. These are taken from
462 the IBC paper. Authors converted the original Fetch environments to a reversible setting by defining
463 a constraint on the block position. The initial and goal distributions are identical to the original Fetch
464 Push and Pick&Place. More details can be found in the IBC paper.

465 **Push (hard):** Different from the original Fetch Push task, in our case walls are added to prevent the
466 block from dropping out of the table. The workspace of the robot arm is also limited. The block
467 is always initialized to a fixed location, and goal distribution during the evaluation is $\mathcal{U}(-0.15, 15)$.
468 Fetch Push used in the IBC paper, the block is limited by joint constraint, which shows unrealistic
469 jittering behaviors near the limits (we observe such phenomenon by running model-based go-explore,
470 the exploration policy prefers to always interact with the block and keep pushing it towards the limit
471 boundary, see videos on our project website¹). Meanwhile, the gripper is blocked, which makes the
472 task easier. In our case, we release the gripper and it can now open and close again which add two
473 more dimension of the state space. We found it is important to release the gripper in our version of
474 Push task, when the block is in corners, it will need to operate the gripper to drag the block escape
475 from corners. The maximum steps the agent can take in 50 during the evaluation. Hard reset will
476 happen after every $1e5$ steps. In the whole training process we performed, it resets 5 times in total.

¹<https://sites.google.com/view/morefreet>

477 **Pick&Place (hard):** We add walls in the same way as we did for Push (hard). We make it more
478 difficult by only evaluating the agent on goals that are in the air. Then it has to learn to perform
479 picking behavior properly, whereas goals on the ground can just be solved by pushing. The goal will
480 be uniformly sampled from a $5 \times 5 \times 10$ cm cubic area above the table. It has the same observation
481 space, action space, initial state and maximum steps with Fetch Push described above. Hard reset will
482 happen after every $1e5$ steps. In the whole training process we performed, it resets 5 times in total.
483 See the visual difference between our Pick&Place and IBC’s in Figure 3. Since the workspace of the
484 robot is limited within the walls as well in Push (hard) and Pick&Place (hard), when the block gets
485 stuck in corners, the robot needs to precisely move to the corner and bring the block back. In contrast,
486 the robot in IBC’s version can move to everywhere, being able to create various circumstance to solve
487 such difficult position.

488 **Ant:** We adapt the AntMaze task from environments² codebase of PEG and change the shape of the
489 maze to square, also change the evaluation goal distribution to be a uniform distribution $\mathcal{U}(2, 3)$ for
490 both x and y location, which lies on the top-left corner of the square. The ant is always initialized to
491 the center point (0, 0) of the square to start from, with uniform noise ($\mathcal{U}(-0.1, 0.1)$) added. The state
492 space is 29D and the action space is 8D. The maximum steps for evaluation is 500. Hard reset will
493 happen after every $2e5$ steps. In the whole training process we performed, it reset 4 times in total.

494 C.2 Baseline Implementations

495 **reset-free PEG:** We extend the official implementation of PEG³ to reset-free setting by 1) removing
496 the reset of environments; 2) optimizing the goal distribution every $H_G + H_E$ steps; 3) keeping all
497 other hyperparameters the same as MoReFree.

498 **IBC:** We use the official implementation from authors⁴ and keep hyperparameters unchanged.

499 **MEDAL:** We follow the official implementation of MEDAL⁵ and use the default setting for experi-
500 ments. Since MEDAL requires demonstrations, for tasks from EARL benchmark, demonstrations
501 are provided. For other environments, we generate demonstrations by executing the final trained
502 MoReFree to collect data. 30 episodes are generated for each task.

503 **R3L:** We implement R3L agent by modifying the FBRL agent from MEDAL codebase. The backward
504 policy is replaced by an exploration policy trained using the random network distillation (RND)
505 objective [2]. The RND implementation we follow is from DI-engine⁶.

506 **Oracle:** This is a episodic SAC agent, we use the implementation from MEDAL codebase and keep
507 all the hyper-parameters unchanged.

508 **MoReFree:** Our agent is built on the model-based go-explore method PEG [14], we extend their
509 codebase by adding back-and-forth goal sampling procedure and training on evaluation initial and
510 goal states in imagination goal-conditioned policy training. See our codebase in the supplemental.

511 C.3 Hyperparameters

512 Train ratio (i.e. Update to Data ratio) is an important hyper-parameter in MBRL. It controls how
513 frequently the agent is trained. Every n steps, a batch of data is sampled from the replay buffer, the
514 world model is trained on the batch, and then policies and value functions are trained in imagination.
515 In all our experiments, we only vary n on different tasks. See the table below for different values on
516 different tasks we used through experiments. MoReFree also introduces a new parameter α , which
517 we keep $\alpha = 0.2$ for all tasks and did not tune it at all. All other hyperparameters we keep the same
518 as the original code base.

519 **Different values for α .** We examine different values of α in MoReFree on Fetch Push task, which
520 affects how much MoReFree focuses on task-relevant goals in exploration and imagination. In
521 Figure 8, we see that introducing a moderate amount of task-relevant goals ($\alpha=0.2$, $\alpha=0.5$) results in

²<https://github.com/edwhu/mrl>

³<https://github.com/penn-pal-lab/peg>

⁴https://github.com/snu-larr/ibc_official

⁵<https://github.com/architsharma97/medal>

⁶https://opendilab.github.io/DI-engine/12_policies/rnd.html

Table 2: Different train ratio we used for different tasks. We keep all other hyperparameters the same as default ones.

PointUMaze	2	Sawyer Door	5	Tabletop	1	Fetch Push	2
Fetch Pick&Place	2	Push (hard)	2	Pick&Place (hard)	2	Ant	2

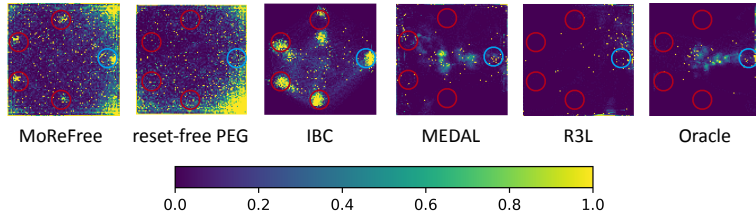


Figure 9: XY state visitation heatmap of the mug in Tabletop of various approaches. MoReFree’s heatmap shows high state diversity while retaining high visitation counts near the task-relevant states (red circles are goal states, the blue circle is the initial state). reset-free PEG also shows diverse exploration, but it over-explores the bottom-right corner which is entirely task-irrelevant. IBC’s bi-directional curriculum leads the exploration shuttles between the initial state and goal states, but fails to explore well. All other methods fail to explore, visited states mostly cluster in few spots.

522 sensible performance, while too many task-relevant goals ($\alpha=0.7$, $\alpha=1.0$) degrades performance. We
 523 use the same value of alpha, 0.2, across all tasks, which showcases MoReFree’s consistency.

524 C.4 Results Clarification

525 In Push and Pick&Place results, we retrieved the final performance of MEDAL directly from the IBC paper (dashed
 526 purple lines) and did not have time to run R3L in these two environments. R3L is shown to be a lot worse than
 527 MEDAL in the MEDAL paper and performs obviously
 528 bad in other tasks shown in Figure 4. In Push (hard) and
 529 Pick&Place (hard), we ran R3L and MEDAL with less
 530 budget since other methods clearly outperform and their
 531 learning curves do not show any evidence for going up.
 532
 533

534 C.5 Resource Usage

535 We submit jobs on a cluster with Nvidia 2080, 3090 and
 536 A100 GPUs. Our model-based experiments take 1-2 days
 537 to finish, and the model-free baselines take half day to one
 538 day to run.

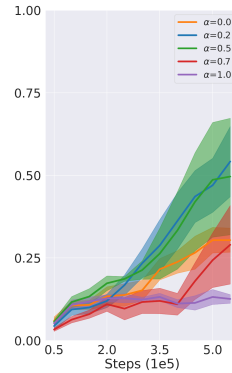


Figure 8: Performance of MoReFree with different values of α in Push (hard).

539 D More Visualizations on Replay Buffer

540 We visualize the replay buffer of different agents on more tasks. See Figure 9 for XY location of
 541 the mug in Tabletop, Figure 11 for XY location data of the agent in PointUMaze, Figure 10 for
 542 XZ location of the block in Pick&Place (hard) and Figure 12 for XY location data of the block in
 543 Push (hard) and Pick&Place (hard). Overall, we see MoReFree explores the whole state space better.
 544 Meanwhile, due to back-and-forth procedure, MoReFree collects more data near initial / goal states,
 545 which are important for the evaluation. However, IBC, MEDAL, R3L and Oracle all fail to explore
 546 well; their heatmaps are mostly populated with low visitation cells.

547 E Detailed Ablations

548 We report learning curves for each variant agent we ablate in Section 5.3 on every task in Figure 13.
 549 Since MoReFree does not learn at all in Sawyer Door task, we exclude the ablation for it. In each

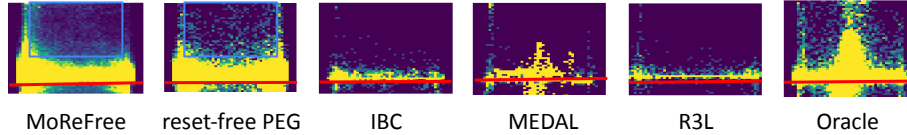


Figure 10: XZ state visitation heatmap of the block in Pick&Place (hard). States above the red line are in the air, which are crucial for solving the picking task. Two MBRL methods collect more data diversely in the air, while other reset-free methods barely pick up the block.

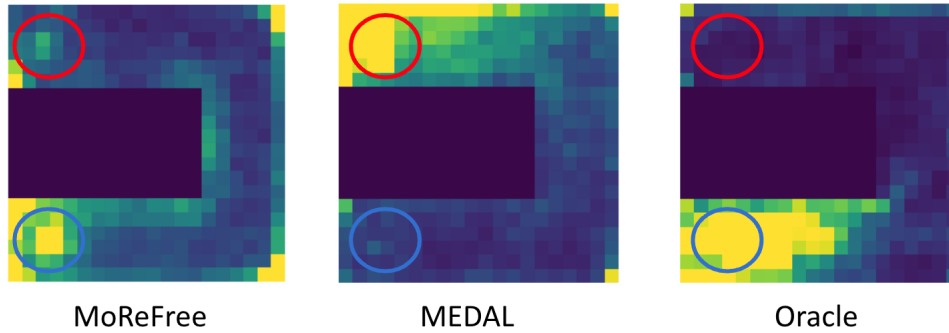


Figure 11: State visitation heatmap on point maze. MoReFree has special focuses on both initial state (blue circles) corner and goal state (red circles), while explore much uniformly. MEDAL collects lots of data near the goal state and little data on the initial state. Both MEDAL and Oracle explore less extensively.

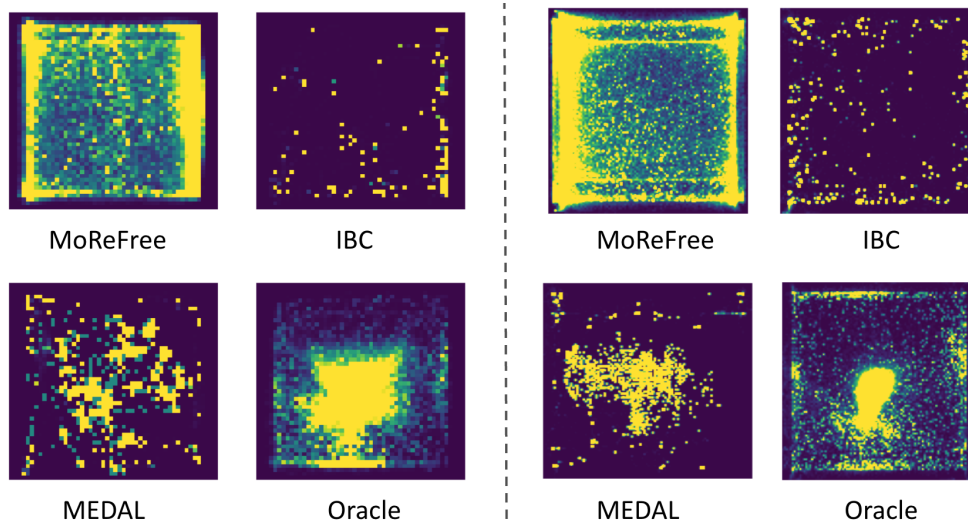


Figure 12: Block state visitation heatmap on Fetch Push (left) and Fetch Pick&Place (right) of different agents. MoReFree better explores the whole state space, while IBC and MEDAL do not have too much interactions with the block, thus lighted areas are scattered everywhere.

550 task, MoReFree is better or on par with all other ablations. Through learning curves, we see different
 551 components contribute differently on different tasks.

552 We further analyze the ablation on PointUMaze as an example by visualizing the replay buffer of
 553 different variants, see Figure 14. In the performance on PointUMaze from Figure 13, sampling
 554 exploratory goals for data collection is important (MF w/o Explore & Imag. outperforms other
 555 ablations). But we see in 14, MF w/o Explore & Imag. does not have focus on the initial / goal state
 556 which we care about for the evaluation, which makes it slightly worse than MoReFree. MF with Only
 557 Task Goals has a strong preference on initial / goal state, we think it is because in the later phase of

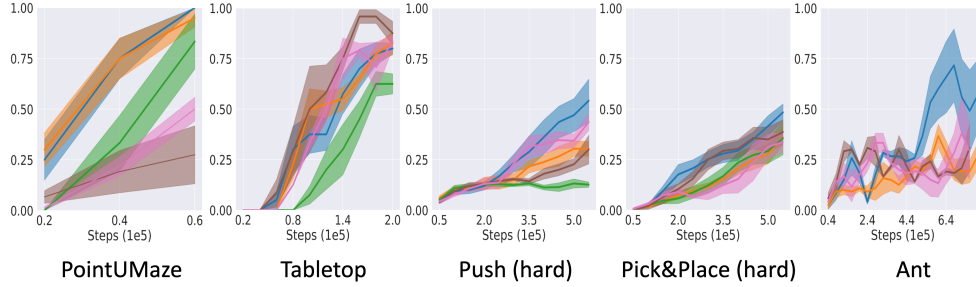


Figure 13: Learning curves of ablation study on 5 tasks. We see different components contribute differently in different tasks. For instance, in Tabletop, **MF w/o Imag.** even performs better than MoReFree, maybe because the whole state space can be explored quickly, then randomly sampling states from the replay buffer as goals for training already has good coverage on evaluation initial / goal states.

558 the training when the agent is able to solve the task, it goes back-and-forth consistently to collect
 559 data. But in the early phase of the training, it might lack exploration which causes the degraded
 560 performance compare with MoReFree. MF w/o Explore and MF w/o Imag. only either go to initial /
 561 goal state for data collection and do not practice on it during the imagination training, or practice
 562 without really going, which both does not form the positive cycle, and end up with poor performance.

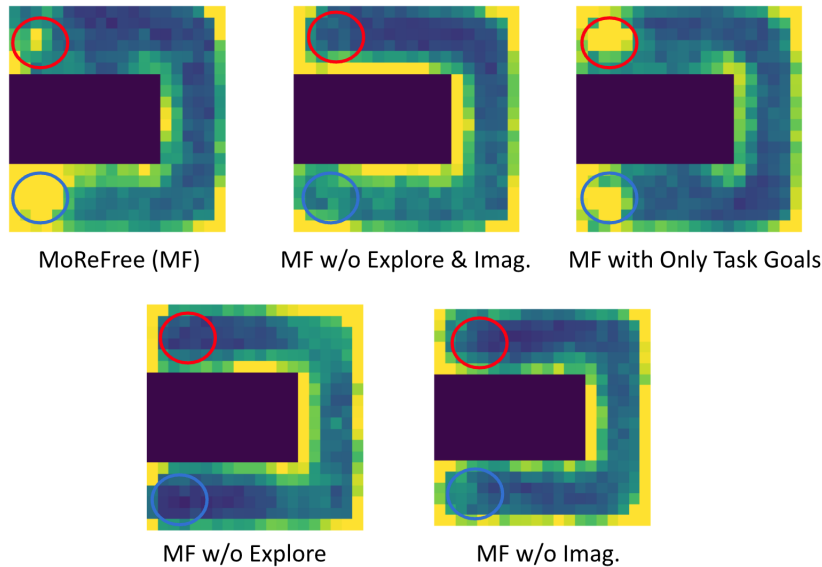


Figure 14: State visitation heatmap on PointUMaze task of all ablations. Red circles are evaluation goal states and blues are initial states. We see MoReFree collect good amount of data near initial / goal states while stronger exploration. MF w/o Explore and MF w/o Imag. could not gather task-relative data, which further causes poor performance.

563

564 F MBRL on Sawyer Door

565 We investigate why two MBRL methods fail on Sawyer Door tasks. Note that MoReFree is able to
 566 solve intermediate goals such as closing the door in some angles, but is unable to solve the original
 567 IBC evaluation goal (see website for more videos).

568 We simplify Sawyer Door task by limiting the movement range of the robot to a box and also having a
 569 block holds the door to prevent it from opening it too much, see Figure 15. Although MBRL methods

570 are trained on the simplified environment, we see learning curves on Sawyer Door are completely
 571 flat in Figure 4, compared with other baselines trained on the original task. We wonder why MBRL
 572 methods can show the same performance and gain benefits as it does in other environments.

573 MoReFree and reset-free PEG use DreamerV2 as backbone agents and extend it to reset-free settings.
 574 We hypothesize that Dreamer itself, even under the episodic setting with task reward function, would
 575 not work well. If that’s the case, then MBRL methods in the reset-free setting with self-supervised
 576 reward function would almost certainly not work either. For example, if the backbone agent cannot
 577 model the dynamics precisely, then policy learning, dynamical distance reward learning, will be
 degraded.

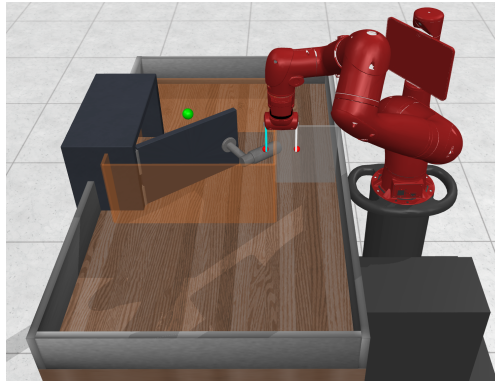


Figure 15: Simplified version of Sawyer Door. Orange walls show the limited workspace for the robot arm, and a grey wall is added to limit the movement of the door. The door can only move to maximum 60 degrees.

578

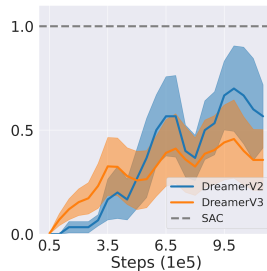


Figure 16: Performance of DreamerV2 and V3 on episodic Sawyer Door task. SAC can solve the task in 200k steps, while after 1 million steps MBRL is still not able to steadily solve the task.

579 We then run the underlying MBRL backbones under the episodic setting. Figure 16 shows DreamerV2⁷, and Dreamerv3⁸ struggle to solve the task, while model-free method SAC can steadily solve the task after 200k steps. This might be a potential reason that MBRL methods do not work on the more difficult reset-free setting. We hypothesize that the combination of the sparse environmental reward and dynamics of the door result in a hard prediction problem for world modelling approaches. We leave further investigation for the future work.

585 G More Analysis on Fetch Environments

586 Although IBC gains good final performance in Push and Pick&Place, it starts learning late compared
 587 with MBRL methods and fails entirely in our harder versions. We suspect IBC might need more computational budget to start learning in harder tasks. Thus we train IBC with two millions environment steps and results in Figure 17 show that it still fails to solve the harder version of Push.

⁷<https://github.com/danijar/dreamerv2>

⁸<https://github.com/danijar/dreamerv3>

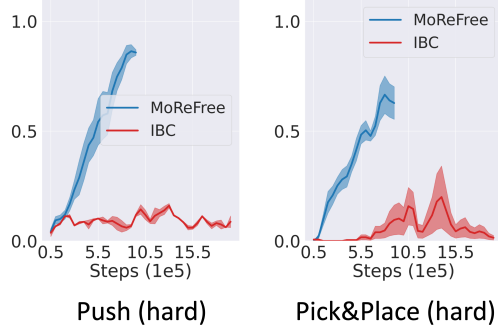


Figure 17: Longer training of IBC in our Fetch tasks, where the state space is larger and artificial constraints are replaced with surrounded walls. IBC still can not learn meaningful behaviors.

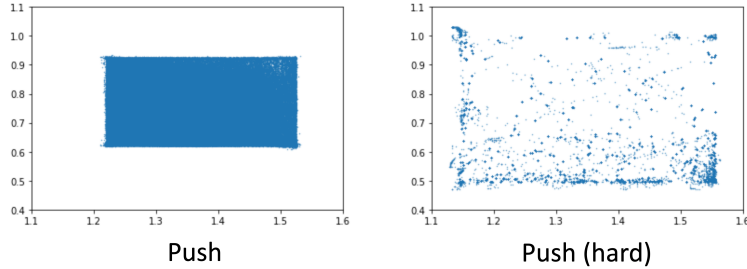


Figure 18: XY location of the block collected by IBC on Push (hard) and its original version (Push). IBC covers the whole state space very well in Push while fails in Push (hard), where the block stays for long time in corners or areas next to walls.

590 Figure 18 shows 600k data of the object (XY view) collected by IBC on our Push (hard) and IBC’s
 591 Push. We see the block stays in corners or next to walls a lot in Push (hard), while goes everywhere
 592 and covers the whole space in IBC’s Push, indicating object interaction is more difficult in Push
 593 (hard) due to the larger state space, surrounded walls and limited work space. In IBC’s Push, the
 594 block can bounce back when it hits the limit of joint constraints. However, in Push (hard), the block
 595 needs to be explicitly brought back from the corner or walls, requiring more sophisticated behaviors.
 596 Meanwhile, larger size of the limited area (our version is 3× larger than IBC’s.) also increases the
 597 difficulty of the task.

598 H Analysis on R3L

599 R3L trains two policies, one for reaching the goal and another that brings the agent to novel states.
 600 The goal-reaching policy is trained using a learned classifier to classify the goal state and other states.
 601 Original R3L takes images as inputs, thus the trained classifier can successfully classify goal images
 602 from random state images. In our work, we use low-dimensional state input. Outputs of the trained
 603 classifier on the whole state space of PointUMaze is shown in Figure 19. We see that the classifier
 604 learns to output higher values for states close to the goal state (red dot) and lower values for states
 605 further away. Nonetheless, due to the smoothness of the output scope, states near the initial state
 606 (blue circle) that are numerically closer but spatially further to the goal state also have higher values.
 607 R3L agent trained using such reward function will always tend to follow states with higher values to
 608 the corner instead of going forward. See the website for more videos. These trained reward functions
 609 are misleading for learning reasonable policies which result in poor performance we see in Figure 4.

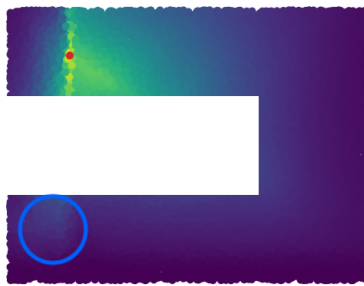


Figure 19: Outputs of the learned classifier on the whole state space. Due to the smoothness of the output scope, states near the initial state (blue circle) also have higher values.