# Using Discrete Overlapping Partitions for Count-Based Exploration

Anonymous authors Paper under double-blind review

Keywords: Exploration, Discrete representation, Intrinsic reward

## Summary

Exploration remains an open problem in reinforcement learning. Ideally, a useful exploration method should efficiently explore sparse reward environments, scale to large environments, and be simple to implement. Most agents use random-based methods like  $\epsilon$ -greedy due to their simplicity and low computational cost. However, these methods struggle in sparse reward settings and take a long time to converge. Count-based methods encourage exploration in less-visited areas but do not scale well to large environments. Extensions to count-based methods to work with function approximation improve performance in complex environments like Montezuma's Revenge but are rarely used due to their computational and implementation complexity. We propose a new method that achieves all three desiderata simultaneously. Our exploration method handles large environments by maintaining counts within multiple overlapping partitions of states to derive exploration bonuses. We evaluate our algorithms on three continuous observation environments where count-based methods cannot be applied including MiniGrid DoorKey with image-based observations.

## **Contribution**(s)

1. We introduce a novel, yet simple, method for exploration that uses overlapping partitions of the state space in order to maintain counts, and then derives an exploration bonus from aggregating these counts to encourage exploration to less-visited states.

**Context:** Prior work has explored extending count-based exploration methods to large environments, most notably pseudo-counts (Bellemare et al., 2016), though these methods are rarely used due to their complexity. In comparison, our method involves obtaining a discrete representation of the environment, maintaining a count for each cell, and computing the intrinsic reward at each timestep. In this paper, we use tile-coding and VQ-VAE to obtain discrete representations of the environment, but other methods could also be used, such as the representations learned in DreamerV2 (Hafner et al., 2020) or fuzzy tiling activation on the last layer of a neural network (Pan et al., 2019).

2. We show that the new method can explore in sparse reward settings and scale to large environments that use complex function approximators. We use continuous river-swim as a known hard exploration environment and the MiniGrid DoorKey environment as an example image-based domain that uses complex function approximation.

**Context:** We build on top of PPO (Schulman et al., 2017) by adding in our exploration bonuses into the optimized return. We compare to the most common approaches to exploration in PPO: random exploration such as epsilon-greedy and entropy regularization on the policy gradient loss. We do not compare to pseudo-count methods as the computational and implementation overhead of these methods make them already very rarely used (Wurman et al., 2022; Bellemare et al., 2020; Berner et al., 2019).

# **Using Discrete Overlapping Partitions for Count-Based Exploration**

#### Anonymous authors

Paper under double-blind review

#### Abstract

1	Exploration remains an open problem in reinforcement learning. Ideally, a useful
2	exploration method should efficiently explore sparse reward environments, scale to
3	large environments, and be simple to implement. Most agents use random-based
4	methods like $\epsilon$ -greedy due to their simplicity and low computational cost. How-
5	ever, these methods struggle in sparse reward settings and take a long time to con-
6	verge. Count-based methods encourage exploration in less-visited areas but do not
7	scale well to large environments. Extensions to count-based methods to work with
8	function approximation improve performance in complex environments like Mon-
9	tezuma's Revenge but are rarely used due to their computational and implementation
10	complexity. We propose a new method that achieves all three desiderata simultane-
11	ously. Our exploration method handles large environments by maintaining counts
12	within multiple overlapping partitions of states to derive exploration bonuses. We
13	evaluate our algorithms on three continuous observation environments where count-
14	based methods cannot be applied including MiniGrid DoorKey with image-based
15	observations.

#### 16 1 Introduction

17 Exploration is a fundamental challenge in reinforcement learning. What experience is gathered 18 is often as important as the process of learning from it or generalization to new experience. One 19 primitive approach is to use random exploration, e.g., a random action is selected with probability  $\epsilon$ 20 at every timestep, or an entropy regularization term is included in the training loss to encourage the 21 policy to remain sufficiently random. These approaches often perform extremely poorly in so-called 22 hard exploration settings, most often characterized by sparse rewards. In such settings, the agent 23 relies entirely on this action randomization to stumble on informative rewards before any of the rest 24 of the learning apparatus can make progress.

25 An early improvement to random exploration has been count-based approaches. These are often 26 motivated under the principle of optimism under uncertainty, where state-action counts can be used 27 to construct exploration bonuses that act to create upper confidence bounds on value estimates (e.g., 28 Strehl & Littman, 2008). These count-based methods, though, are not directly applicable in very 29 large or continuous state and action spaces, where an agent may never revisit a state during typical 30 lengths of trials. In such cases, alternative methods that are friendlier to function approximation 31 have been proposed. For example, Density-based pseudo-counts (Bellemare et al., 2016) construct 32 an exploration bonus from a density model learned alongside a value function. Random network dis-33 tillation (Burda et al., 2018) creates a bonus based on the error of predicting the output of randomly 34 initialized fixed network of the state input — the prediction improves with observations, decay-35 ing the bonus like counts. It would seem we have an ideal answer that is function approximation 36 friendly.

However, celebrated results such as GT-Sophy (Wurman et al., 2022) outperforming professional 37 38 Gran Turismo drivers or navigation of stratospheric balloons (Bellemare et al., 2020) use simplis-39 tic random exploration (entropy regularization and random-walk policies, respectively). Further-40 more, Dreamer v3 (Hafner et al., 2023), a general-purpose algorithm achieving state-of-the-art results across many tasks, also eschews directed exploration for simple random exploration from 41 42 entropy regularization. We propose that another consideration is the simplicity of the exploration 43 method. Pseudo-counts and random-network distillation require significant computational overhead 44 with training a wholly separate network (or two) in addition training of a value or policy network. 45 They also are complicated to implement and correspondingly complicated to make work. Dabney 46 et al. (2020) argues that the complexity of many exploration methods may even limit their generality, 47 suggesting that is the reason  $\epsilon$ -greedy remains so common.

In summary, there are three desiderata for exploration: effective directed exploration especially in sparse reward settings, function approximation friendliness, and simplicity in implementation, computation, and applicability. In this work we propose a directed exploration approach nearly as simple as count-based methods for tabular function approximation while remaining friendly to common large-scale function approximation schemes, achieving all three desiderate.

53 Our approach exploits two things. First, many modern function approximation schemes are observ-54 ing considerable benefit in constructing discrete, but combinatorial, representations. For example, 55 Dreamer (Hafner et al., 2020; 2023) builds discrete states for its world model made up of 32 variables 56 each able to take on 32 discrete values. Each of these easily countable. VQ-VAE (Van Den Oord 57 et al., 2017) gives down-scaled discrete representations of an image, where each latent pixel takes 58 on finite discrete index values; again, easily countable. Fuzzy tiling activations (Pan et al., 2019) 59 provide a discretized final layer representation, aimed at making learning more stable to covariate 60 shift, but, also easily countable. These discrete combinatorial representations, can be thought of as 61 providing multiple overlapping partitions of the state space: each variable partitions the states by the values it can take. These will be what we count. Second, we aggregate the counts associated 62 with each variable to build a single associated count for a state and then use this as a traditional 63 exploration bonus. 64

We evaluate this approach on simple but challenging toy exploration problems with continuous state spaces and finally on a MiniGrid (Chevalier-Boisvert et al. 2024) sparse reward task using image-

spaces and finally on a MiniGrid (Chevalier-Boisvert et al., 2024) sparse reward task using image-

based observations. We show that our *simple* count-based method is both *function approximation* 

68 *friendly* and *effective at directing exploration* in sparse reward settings.

#### 69 2 Background

We begin by introducing the reinforcement learning (RL) problem setting and discuss two popular solution methods, Q-learning (Watkins & Dayan, 1992) and Proximal Policy Optimization (PPO; Schulman et al., 2017), which we use alongside our exploration method. We then present the common exploration strategies associated with these methods, specifically  $\epsilon$ -greedy for Q-learning and entropy regularization for PPO. We also describe two discretized state representation techniques, tile-coding and VQ-VAE, which are used in Section 4 to obtain overlapping partitions for our exploration method.

77 Consider a Markov Decision Process (MDP), defined as a tuple  $\langle S, A, R, P, \gamma \rangle$  where S represents the state space,  $\mathcal{A}$  represents the action space, R(s, a) is the reward function giving the expected 78 79 reward for taking action a from state s, P(s'|s, a) is the transition function giving the probability of transitioning from state s to state s' given action a, and  $\gamma$  is the discount factor. In this paper we 80 81 focus on model-free RL where the transition function P is unknown and is not learned by the agent. 82 The environment can be finite and discrete, where  $|\mathcal{S}| = d_{\mathcal{S}}$  and  $|\mathcal{A}| = d_{\mathcal{A}}$ , or it can be continuous, where  $S \subseteq \mathbb{R}^{d_S}$  and  $A \subseteq \mathbb{R}^{d_A}$ . When we define counts on state-action pairs, we focus on the case 83 84 where the state and action space is finite and discrete. When we introduce our new approach we 85 will allow the state space to either be continuous or discrete (and large), but the action space for 86 simplicity will remain discrete.

The goal of an RL agent is to take actions to maximize the discounted sum of rewards,  $\sum_{t=1}^{\infty} \gamma^{t-1} r_t$ , 87 it receives over time where  $\gamma \in [0,1)$  (Sutton, 2018). We can represent the agent's choice of 88 action with a policy  $\pi(s)$  giving a probability distribution over  $\mathcal{A}$  for state s. To quantify the the 89 90

performance of policy  $\pi$ , we define its state-action value or Q-value as:

$$Q^{\pi}(s,a) \equiv \mathbb{E}^{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r_{t+1} \middle| s_0 = s, a_0 = a \right], \tag{1}$$

where  $\mathbb{E}^{\pi}$  denotes the expected value over all possible trajectories following  $\pi$ . A common goal for 91 RL is to learn an optimal policy  $\pi^* = \arg \max_{\pi} Q^{\pi}(s, a)$ . 92

#### 2.1 Q-learning 93

One common approach to RL is to estimate a value function. The Q-value, can be calculated using 94 95 the recursive Bellman equation

$$Q^{\pi}(s,a) = R(s,a) + \gamma \sum_{s',a'} P(s'|s,a)\pi(a'|s')Q^{\pi}(s',a'),$$
(2)

which is the basis for many learning algorithms, including Q-learning. Q-learning is a temporal 96 97 difference (TD) control algorithm where Q is updated by computing the TD error, the difference in 98 the Bellman equation given a sample of the next state s' and assuming a' is being selected greedily

99 with respect to the current Q estimate. Using  $\alpha$  as a step size, the Q estimate is updated at time t by,

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t [r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)].$$
(3)

Here, we assume we have a finite action and state space, and the Q estimates are stored as a table. 100 101 In order to guarantee convergence to the optimal policy, all states and actions need to be visited infinitely often, i.e., we need an effective exploration method. This is often achieved through selecting 102 103 actions  $\epsilon$ -greedy with respect to the current Q estimates.

#### 104 2.2 Proximal Policy Optimization (PPO)

105 Another common approach is to represent and learn a policy directly. Let  $\pi_{\theta}(a|s)$  be a parameter-106 ized representation of the policy with parameters  $\theta$ , e.g., using a neural network. The goal becomes to directly update  $\theta$  to take actions that accumulate more reward, typically through gradient up-107 108 dates. One such approach is PPO (Schulman et al., 2017), which also learns a value function used 109 to reduce variance in the policy updates. PPO collects experience through interactions with the en-110 vironment and optimizes an objective function using minibatch updates from multiple trajectories. 111 PPO's objective to maximize can be written as:

$$J_t(\theta) = \mathbb{E}_t \left[ J_t^{\text{CLIP}}(\theta) - c_1 L_t^{\text{VF}}(\theta) + c_2 H[\pi_\theta] \right], \tag{4}$$

where  $c_1, c_2$  are coefficients,  $J_t^{\text{CLIP}}$  is the policy improvement objective (clipped to prevent the policy from changing too drastically on each update, i.e., keep it "proximal"),  $L_t^{\text{VF}}$  is the value function 112 113 error and  $H[\pi_{\theta}]$  is an entropy regularizer. Our main focus is the entropy regularizer, which serves as 114 PPO's main exploration mechanism. 115

#### 2.3 Exploration Bonuses 116

117 An influential approach to exploration is the principle of Optimism in the Face of Uncertainty (OFU). 118 Instead of greedily choosing actions based solely on estimated rewards (with or without undirected 119 randomness), the agent assumes that less visited states and actions may have a higher value than 120 estimated. This is typically done with some form of exploration bonus added to the environment re-121 ward to drive a greedy policy to seek less visited states and actions. As the agent explores and visits more states and actions, it reduces its uncertainty and the associated exploration bonus. The Upper Confidence Bound (UCB; Auer, 2002) algorithm belongs to the OFU family. It was introduced in the multi-armed bandit setting, but has also been used extensively in the MDP setting (e.g., ?Bellemare et al., 2016). The agent operates as if the observed reward is  $r_t = R(s_t, a_t) + \beta R^+(s_t, a_t)$ where  $R^+(s_t, a_t)$  denotes the exploration bonus, and  $\beta$  is the coefficient balancing exploration and exploitation. UCB suggests the following count-based exploration bonus:

$$R^{+}(s,a) = \sqrt{\frac{2\log N(s)}{N(s,a)}}$$
(5)

where N(s, a) is the state-action pair visitation count and  $N(s) = \sum_{a} N(s, a)$ . Maintaining counts for every state-action pair is not applicable in large or continuous environments. Not only is it intractable to store such counts, but an agent may never revisit a state during typical lengths of trials, and so will not effectively direct exploration. Density-based pseudo-counts (Bellemare et al., 2016) and random network distillation (Burda et al., 2018) are attempts to devise count-like alternatives for Equation 5 that are function approximation friendly. However, these are complex in both implementation and computation, and so do not meet the stated desiderata.

#### 135 2.4 State Discretization

136 Our approach makes use of a growing trend in RL: learned discretized state representations. In 137 order to have high expressivity, these representations are typically combinatoric, i.e., the state is 138 represented as the composition of a set of discrete variables. Each discrete variable of a state,  $v_i(s) \in$  $[n] \equiv \{1, \dots, n_i\}$ , can be thought of as partitioning states into a finite number of cells indexed by 139 140 positive integers. And the overall representation  $\phi(s) = v_1(s) \times \ldots \times v_n(s)$  is the Cartesian product 141 of these multiple overlapping discrete partitions. Such a combinatoric representation is not new: 142 CMACs (Albus, 1975) and tile-coding (Sutton, 1995) are examples of such representations. In tile-143 coding, a grid or tiling is defined on top of a vector representation of the state so as to partition the 144 states into tiles. By offsetting the tiling by different amounts, multiple overlapping partitions can 145 be created, where the state is represented by the set of discrete tiles that contains it in each tiling. 146 Large tiles then give a sense of generalization, as states that are distant will share some tiles in the 147 representation, while a large number of tilings gives expressivity as nearby states can still be in 148 different tiles for some tiling.

The recent trend is to learn such discrete representations. Dreamer (Hafner et al., 2020; 2023) is one 149 150 such example where a state's learned latent representation is a binary matrix with 32 rows of one-hot 151 encoded vectors of length 32. Each row, like tile-coding, partitions the state-space into 32 cells, with 152 the one-hot encoding identifying the cell. Expressivity comes from combinatorics of representing 153 the state with 32 such discrete partitions. Vector Quantized Variational Auto Encoders (VQ-VAE; 154 Van Den Oord et al., 2017) is another example, and one we will use in our experiments. A VQ-VAE 155 consists of an encoder, a quantization step, and a decoder, aimed at finding a latent representation of 156 images that ideally respects an image's compositionality. By passing the image input to the encoder we can learn a representation of the input  $z_e(x) \in \mathbb{R}^{n \times m}$ , which gives n vectors of length m. These 157 vectors are compared to a dictionary of embedding vectors  $e \in \mathbb{R}^{k \times m}$ . For vector *i*, let  $v_i \in [k]$  be 158 the row index of the chosen vector of e. These chosen vectors,  $z_q \in \mathbb{R}^{n \times m}$ , are passed through a 159 decoder to reconstruct the given input. The encoder, decoder, and dictionary of embedding vectors 160 161 are all learned in VQ-VAE. Notice that the Cartesian product  $v_1(s) \times \ldots \times v_n(s)$  gives multiple 162 overlapping partitions of its input, just as with tile-coding. Locally sensitive hashing (Indyk & 163 Motwani, 1998; Bellemare et al., 2013) and fuzzy tiling activations (Pan et al., 2019) are two other 164 examples of fixed and learned discrete representations, respectively.

#### 165 **3 Method**

166 Consider the Continuous Riverswim (Pan et al., 2018) environment (see Figure 15), which expands

167 the discrete 6-state classic environment, to the [0,1] interval. The environment has two actions:



Figure 1: Continuous Riverswim illustration. The red and green regions indicate low- and high-reward areas, respectively. The black dot represents the agent within three overlapping tilings.

168 LEFT and RIGHT. LEFT moves the agent 0.1 units to the left, while RIGHT moves the agent 0.1169 units right (with 30% probability), 0.1 units left (with 10% probability), or keeps it in place (with 170 70% probability). After moving, the agent's position is perturbed by a sample from a zero-mean 171 Gaussian with standard deviation 0.01. If the agent takes RIGHT from a state > 0.95 it receives a 172 reward of 10,000; LEFT from a state < 0.05 a reward of 5, otherwise 0;  $\gamma = 0.99$ . The optimal 173 policy is to always choose RIGHT, however it is difficult and unlikely to observe the large reward 174 doing so, while the small reward for LEFT is easily found. While count-based methods efficiently 175 explore in the discrete version of Riverswim, it cannot be applied in Continuous Riverswim.

176 Now consider applying tile-coding to the Continuous Riverswim state space, where each partition 177 is just a one-dimensional tiling (see the tiling illustration in Figure 1), and each subsequent tiling 178 is offset to create different boundaries between cells in each partition. A single state can then be 179 represented by the cells from each partition that contains that state. It is these cells, i.e., tiles, that 180 our approach will count, so each state visited leads to incrementing the count of all n cells containing 181 the state. To construct an exploration bonus from the n counts, we use an aggregation function to 182 turn a set of counts into a single count and then apply Equation 5 to the aggregate count. Possibilities 183 for the aggregation function include average, max, and min. The choice of min for the aggregation 184 is suggestive of the count-min-sketch (Cormode & Muthukrishnan, 2005) for estimating element 185 counts with sublinear memory, which also counts cells of multiple overlapping partitions but uses 186 perfect pairwise-independent hash functions as the partitions. The choice of average is suggestive 187 of tile-coding itself, where a value function is represented by the average value of each cell.

Formally, imagine a set of overlapping partitions  $\{P_i\}_{i=1...n}$  where each partition  $P_i : S \to [k]$  maps to k possible cells. We can represent the set of all cells with  $C = \{(i, j) \mid i \in [n], j \in [k]\}$ . We keep counts on every cell action pair,  $N(c \in C, a \in A)$ . Given an aggregation function  $A : \mathbb{N}^n \to \mathbb{R}$ we construct a state action count as  $N(s, a) \equiv A(\{N(P_i(s), a)\}_{i=1...n})$ , which then gets used to create a UCB exploration bonus with Equation 5.

193 Tile-coding becomes impractical in high-dimensional spaces because the number of required tilings 194 increases exponentially with the number of dimensions to maintain a fixed resolution. Furthermore, 195 as a fixed representation it does not adapt its expressivity to relevant differences in states. Hence, 196 it would not be practical, for example, for image-based state representations. However, the general 197 process can be applied to any discrete representation with multiple overlapping partitions, including 198 the recent trends of learned discrete representations, such as those discussed in Section 2.4. In 199 particular, we will experiment with exploring using counts on the discrete latent spaces produced by 200 VO-VAE.

#### 201 4 Results

202 We evaluate our method across three distinct environments with different characteristics to analyze 203 its effectiveness in various reinforcement learning settings. First, we use the Continuous River-204 Swim (Pan et al., 2018) environment, a small continuous environment with a distractor reward that 205 makes it a challenging exploration problem for Q-learning. Second, we evaluate on the Mountain 206 Car (Moore, 1990) environment, a classic control task with sparse rewards, using PPO. Lastly, we 207 evaluate our method on the MiniGrid DoorKey (Chevalier-Boisvert et al., 2024) environment, where 208 the agent receives image-based observations using VQ-VAE representations. We conducted exper-209 iments on both the standard DoorKey setting and a modified version with shorter episodes to make 210 the exploration problem harder. Unless otherwise noted, all results are averages over 50 runs and 211 shaded regions show 95% confidence intervals. The exact details of all of the environments as well 212 as the hyperparameters for the experiments are all reported in the supplementary material.

#### 213 4.1 Continuous Riverswim

214 Continuous Riverswim is a challenging exploration task 215 that requires a directed exploration strategy. The diffi-216 culty of a random walk to reach the large reward and 217 the presence of the small distractor reward increases the 218 difficulty of exploration, as the agent prematurely may 219 seek the suboptimal reward. We trained agents using 220 Q-learning with different exploration methods:  $\epsilon$ -greedy, 221 state aggregation (discrete representation with no over-222 lapping partitions), and multiple overlapping partitions 223 using tile-coding with two different aggregation func-224 tions: average and min. The same discrete representation 225 was used for both learning and exploration bonuses, with 226 the number of tilings and tiles not optimized for either 227 exploration or learning. Figure 2 shows the total undis-228 counted return over 100,000 steps of learning. We see that 229 using overlapping discrete representations with either of 230 the aggregation functions yields a higher return than state 231 aggregation and  $\epsilon$ -greedy, as it more reliably found the 232 action leading to the large reward. Amongst the two ag-233 gregation functions we see that min is a slower to begin



Figure 2: Continuous RiverSwim showing total undiscounted return over 100,000 timesteps for Q-learning with different exploration methods. "Tilecoding(X)" denotes our proposed approach with different aggregation functions. State Aggregation denotes using counts with a single partition of the state space.

exploiting than average but more reliably exploits the large reward. Note that despite tuning  $\epsilon$  in  $\epsilon$ greedy, testing values  $\epsilon \in \{0.1, 0.2, 0.5\}$ , none of the runs discovered the large reward, consistently exploiting the suboptimal reward.

#### 237 4.2 Mountain Car

238 Mountain Car is a classic control environment with con-239 tinuous state and sparse reward (rewards are -1 except at 240 episode termination). In this environment, we use PPO 241 for learning as it is known to struggle in Mountain Car, 242 while using tile-coding to obtain discrete state represen-243 tations for exploration bonuses. Similar to Continuous 244 Riverswim (Section 4.1) we do not tune the number of 245 tilings or tiles. For PPO's policy and value networks, 246 we use an MLP with position and velocity as input and 247 two hidden layers, each containing 64 units, with tanh as 248 the activation function. Figure 3 presents the results for 249 entropy regularization and our count-based method using



Figure 3: Mountain Car showing episodic return over 500k timesteps for PPO using tile-coding counts and vanilla PPO with entropy regularization.

average as the aggregation function. As shown in Figure 3, PPO with our count-based exploration learns approximately twice as fast as PPO with entropy regularization. Count-based exploration also demonstrates lower variance in the episodic return, as indicated by the shaded confidence intervals. This is likely due to PPO with entropy regularization failing to reliably learn a policy to reach the goal within the training duration on some trials, which demonstrates the weakness of relying on strictly random exploration rather than directed exploration to reduce uncertainty as with exploration bonuses.

#### 257 4.3 MiniGrid Doorkey

258 Finally, we test our approach in the MiniGrid DoorKey environment (Chevalier-Boisvert et al., 259 2024). MiniGrid DoorKey at its hear is a gridworld consisting of two separate rooms divided by 260 a wall and a locked door. The agent's goal is to pick up the key from one room, unlock the door, and 261 reach the green tile in the second room. However, the agent receives an image of the environment as 262 its state observation. The available actions include turning right, turning left, moving forward, pick-263 ing up an object, and using an object (e.g. using the key to unlock the door). The agent receives a reward of 0 on each timestep, and a reward of  $1-0.9\frac{t}{T}$  upon reaching the goal, where t is the current 264 265 timestep and T is the maximum number of timesteps in an episode. The environment is stochas-266 tic: each time the agent selects an action, there is a 10% probability that a different random action 267 will be executed instead. Figure 4a depicts the 8×8 MiniGrid DoorKey environment used in our 268 experiments, illustrating the agent's starting position, goal location, and the key's location. Finally, 269 an episode terminates after 1,000 timesteps if the goal has not been reached. We observed that this 270 makes for an easy exploration policy as a 1,000 step random-walk has a significant chance of reach 271 the goal, which is under 20 steps from anywhere in the environment. We also experimented with 272 a more challenging variant, where episodes were terminated after 100 timesteps, making directed 273 exploration more important.

To encode the environment into multiple overlapping partitions, we train a VQ-VAE model. These representations are later used to train a PPO agent and maintain state visitation counts. We follow the model-free experimental settings of Meyer et al. (2023). A VQ-VAE model is pre-trained with 500k steps of data gathered from a random walk in the environment. The discrete representation was then used as both the input state for PPO's policy and value function, as well as the basis for

279 defining exploration bonuses from the cell counts.



Figure 4: Minigrid DoorKey showing the environment (a); episodic return with 1,000 step episodes (b) and 100 step episodes (c). Results averaged over 30 runs. Training shown starting after 500k steps of pre-training the VQ-VAE representation.

Figure 4b presents our results on the standard MiniGrid DoorKey environment, where the maximum number of timesteps per episode is 1,000. Agents were trained for 200k steps, which we show starting from 500k to account for the VQ-VAE pre-training. We compare a PPO agent using entropy regularization for exploration with two PPO agents employing count-based exploration, each with different exploration coefficients  $\beta$ . Our approach with  $\beta = 0.001$  and entropy regularization exhibit 285 similar performance and converge to an optimal policy within 100K from the start of PPO training. 286 However,  $\beta = 0.01$  is still learning after 200k steps, which we believe is due to excessive exploration 287 caused by the large exploration bonus coefficient. We believe the agent must spend significantly 288 more time visiting states throughout the environment to sufficiently collapse its uncertainty before 289 the environment reward dominates. State distribution plots in Supplementary Materials S1 confirm 290 this behavior: even after solving the task, the agent with  $\beta = 0.01$  continues to explore other regions, 291 reducing overall efficiency. These results show the importance of balancing between exploring 292 enough to identify the optimal policy and not over-exploring, which can delay exploiting with one's 293 learned policy.

294 What about the harder challenge of shorter episodes, increasing the need to efficiently explore before 295 restarting in the initial state. Figure 4c shows the average episodic return for the same agents. PPO 296 with entropy regularization regularly failed to reach the goal even once in 200k steps, exhibiting no 297 learning at all, as it relied only on a random walk to observe non-zero return. Our approach showed 298 that it could reliably learn good policies in this hard exploration task, with the larger exploration 299 bonus ( $\beta = 0.01$ ) proving to learn slightly faster, again emphasizing the need to balance exploration 300 with the difficulty of the task. This improved exploration was achieved by simple counts on discrete 301 quantities, requiring almost no additional computation and or complexity of implementation.

#### 302 5 Conclusion

303 We introduced a simple count-based method for exploration that exploits discrete representations 304 consisting of multiple overlapping partitions of the state space. We show that the method can do 305 efficient exploration in hard exploration problems, such as Continuous Riverswim, as well as being 306 function approximation friendly while using VQ-VAE representations of image-based observations. 307 This method addresses the three desiderata of simple, effective and function approximation friendly. 308 There are a number of future directions that would be valuable to explore. First, a thorough inves-309 tigation of aggregation functions is needed. Can some theory of count-min-sketch be extended to 310 this setting where representations are not pairwise-independent hashes? Are there situations where 311 min or average perform better? Might soft-min give additional control over the manner of explo-312 ration? Furthermore, how does this approach fair when the discrete representation is being learned 313 simultaneously? Ultimately, we see this approach as sufficiently simple and effective to replace the 314 surprising resilience of  $\epsilon$ -greedy.

#### 315 **References**

- 316 James S Albus. A new approach to manipulator control: The cerebellar model articulation controller
- 317 (cmac). Journal of Dynamic Systems, Measurement, and Control, 97(3):220–227, 1975.
- 318 P Auer. Finite-time analysis of the multiarmed bandit problem, 2002.
- Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos.
   Unifying count-based exploration and intrinsic motivation. *Advances in neural information pro-*
- 321 *cessing systems*, 29, 2016.
- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environ ment: An evaluation platform for general agents. *Journal of artificial intelligence research*, 47:
   253–279, 2013.
- Marc G Bellemare, Salvatore Candido, Pablo Samuel Castro, Jun Gong, Marlos C Machado, Sub hodeep Moitra, Sameera S Ponda, and Ziyu Wang. Autonomous navigation of stratospheric bal loons using reinforcement learning. *Nature*, 588(7836):77–82, 2020.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy
   Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large
   scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.

- Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network
   distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- 333 Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo Perez-Vicente, Lucas Willems,
- Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Mod ular & customizable reinforcement learning environments for goal-oriented tasks. *Advances in*
- 336 Neural Information Processing Systems, 36, 2024.
- Graham Cormode and Shan Muthukrishnan. An improved data stream summary: the count-min
   sketch and its applications. *Journal of Algorithms*, 55(1):58–75, 2005.
- 339 Will Dabney, Georg Ostrovski, and André Barreto. Temporally-extended  $\epsilon$ -greedy exploration. 340 *arXiv preprint arXiv:2006.01782*, 2020.
- Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with dis crete world models. *arXiv preprint arXiv:2010.02193*, 2020.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains
   through world models, 2023. URL https://arxiv. org/abs/2301.04104, 2023.
- Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of
  dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*,
  pp. 604–613, 1998.
- Edan Meyer, Adam White, and Marlos C Machado. Harnessing discrete representations for continual reinforcement learning. *arXiv preprint arXiv:2312.01203*, 2023.
- Andrew William Moore. Efficient memory-based learning for robot control. Technical report, Uni versity of Cambridge, Computer Laboratory, 1990.
- Yangchen Pan, Muhammad Zaheer, Adam White, Andrew Patterson, and Martha White. Organizing
   experience: a deeper look at replay mechanisms for sample-based planning in continuous state
   domains. *arXiv preprint arXiv:1806.04624*, 2018.
- Yangchen Pan, Kirby Banman, and Martha White. Fuzzy tiling activations: A simple approach to
   learning sparse representations online. *arXiv preprint arXiv:1911.08068*, 2019.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Alexander L Strehl and Michael L Littman. An analysis of model-based interval estimation for
   markov decision processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.
- Richard S Sutton. Generalization in reinforcement learning: Successful examples using sparse
   coarse coding. Advances in neural information processing systems, 8, 1995.
- 363 Richard S Sutton. Reinforcement learning: An introduction. A Bradford Book, 2018.
- Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. Advances in
   *neural information processing systems*, 30, 2017.
- 366 Christopher JCH Watkins and Peter Dayan. Q-learning. Machine learning, 8:279–292, 1992.
- Peter R Wurman, Samuel Barrett, Kenta Kawamoto, James MacGlashan, Kaushik Subramanian,
  Thomas J Walsh, Roberto Capobianco, Alisa Devlic, Franziska Eckert, Florian Fuchs, et al. Outracing champion gran turismo drivers with deep reinforcement learning. *Nature*, 602(7896):223–
  228, 2022.

# Supplementary Materials

371 372 373

The following content was not necessarily subject to peer review.

#### 374 S1 State distribution

In this section, we present the state distribution of three PPO agents using count-based exploration with  $\beta \in \{0.01, 0.001\}$  and PPO with entropy regularization on MiniGrid DoorKey in two variations: the standard version (1,000 maximum steps per episode) and a modified version with a different step limit (1,000 maximum steps per episode). For both environments, we plot the state distribution to visualize how agents explore the state space over time, analyzing visitation frequency and coverage. Additionally, we remove the environmental reward to isolate the effect of intrinsic exploration, allowing us to examine how agents navigate without external incentives.

#### 382 S2 Hyperparameters

383 For Continuous Riverswim we swept over  $\beta \in \{1, 10, 100, 1000, 10000\}, \epsilon \in \{01., 0.2, 0.3, 0.4, 0.5\}$ , and  $\alpha \in \{0.001, 0.1, 1, 10\}$  for all three methods ( $\epsilon$ -greedy, state ag-385 gregation and tile-coding). Before running the main experiment we tuned our hyperparameters and 386 selected parameters that yielded highest average reward.

In the Mountain Car we sweep over  $\alpha \in \{2.5 \times 10^{-3}, 2.5 \times 10^{-4}, 2.5 \times 10^{-5}\}$ , the hyperparameters

were chosen as standard parameters in classic control problem. The hyperparameters for PPO are shown in Table 1.

Table 1:	: Hyper	parameters	for	Mountain	Car
----------	---------	------------	-----	----------	-----

Hyperparameter	Value		
PPO mini batch size	128		
PPO iterations	4		
PPO clipping factor ( $\epsilon$ )	0.2		
PPO value coefficient $(c_1)$	0.5		
PPO entropy coefficient $(c_2)$	0.001		

We use the same hyperparameters as those in Meyer et al. (2023) for PPO and VQ-VAE. We perform a grid search over the step size,  $\alpha \in \{3 \times 10^{-3}, 3 \times 10^{-4}, 3 \times 10^{-5}\}$ , and  $\beta$  coefficient for the count-based exploration algorithms,  $\beta \in \{0.01, 0.001, 0.0001\}$ , keeping all other parameters unchanged. For count-based exploration, the *min* function was chosen as the aggregate function, and the best step size for each coefficient was selected. For every coefficient in both variations of the environment  $\alpha = 3 \times 10^{-5}$  had the highest performance except  $\beta = 0.001$  in the standard version where  $\alpha = 3 \times 10^{-4}$  was chosen.

397 The fixed hyperparameters for PPO and VQ-VAE are listed in Table 2.

#### 398 S3 Environments

#### 399 S3.1 Continuous Riverswim

400 Continuous Riverswim, introduced by (Pan et al., 2018), is a modified version of Riverswim (Strehl 401 & Littman, 2008). On the leftmost side of the environment, a distractor reward encourages the 402 agent to remain there unless sufficient exploration is done, posing a significant challenge. The state 403 space is  $s \in [0, 1]$ , and the actions are  $a \in \{left, right\}$ . Once the direction of movement is

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

beta = 0.01 door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.01	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00		0.00	0.00	0.00	0.00
0.00	0.00	0.00			0.00	0.00	0.00
0.00	0.00	0.00		0.00	0.00	0.00	0.00
0.00		0.01		0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### beta = 0.001 door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.03			0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.02	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.03			0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00		0.00		0.00	0.00	0.00	0.00
0.00		0.01		0.00	0.00	0.00	0.00
0.00		0.01			0.01	0.01	0.00
0.00		0.01		0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door locked, key dropped

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00			0.01	0.00	0.00	0.00	0.00
0.00		0.01	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.04	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

door locked, key dropped

0.00

0.00

0.00

.02 0.01

0.01 0.01 0.01 0.00 0.00 0.00 0.00

0.01 0.01 0.01 0.00 0.00 0.00 0.00

0.00 0.00 0.00 0.00 0.00 0.00

0.00 0.00 0.00 0.00

0.00 0.00

#### PPO door locked, key carry



#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00		0.01		0.00			0.00
0.00		0.00		0.00	0.00	0.00	0.00
0.00		0.00		0.00	0.00	0.00	0.00
0.00							0.00
0.00		0.00		0.00	0.01	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### Figure 5: State distribution for standard DoorKey (1000 steps) after 20k steps of training PPO

¢	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
¢	0.00				0.00	0.00	0.00	0.00
¢	0.00			0.00	0.00	0.00	0.00	0.00
¢	0.00			0.02	0.00	0.00	0.00	0.00
¢	0.00				0.00	0.00	0.00	0.00
c	0.00				0.00	0.00	0.00	0.00
0	0.00			0.04	0.00	0.00	0.00	0.00
0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

door locked, key dropped

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00		0.01	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.04	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door locked, key dropped

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00			0.01	0.00	0.00	0.00	0.00
0.00		0.01	0.00	0.00	0.00	0.00	0.00
0.00			0.02	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.03		0.05	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

beta = 0.01 door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.04	0.03	0.05	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.02	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

beta = 0.001 door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.03	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

PPO door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.04	0.00	0.00	0.00	0.00
0.00			0.03	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.03	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.01	0.01	0.00
0.00	0.01	0.00	0.01	0.00	0.00	0.00	0.00
0.00	0.00	0.00		0.00	0.00	0.00	0.00
0.00		0.01			0.01		0.00
0.00	0.00	0.00		0.00	0.00	0.00	0.00
0.00		0.01		0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00		0.01		0.00			0.00
0.00		0.00		0.00	0.00	0.00	0.00
0.00		0.01		0.00	0.01	0.01	0.00
0.00		0.01					0.00
0.00		0.00		0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00		0.01		0.00		0.01	0.00
0.00	0.01	0.00		0.00	0.00	0.00	0.00
0.00		0.00		0.00	0.00	0.00	0.00
0.00		0.01		0.02			0.00
0.00		0.00	0.01	0.00	0.01	0.01	0.00
0.00				0.00		0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Figure 6: State distribution for standard DoorKey (1000 steps) after 50k steps of training PPO

#### 12

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00		0.01	0.01	0.00	0.00	0.00	0.00
0.00		0.00	0.00	0.00	0.00	0.00	0.00
0.00			0.04	0.00	0.00	0.00	0.00
0.00			0.05	0.00	0.00	0.00	0.00
0.00			0.05	0.00	0.00	0.00	0.00
0.00			0.04	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

## door locked, key dropped

	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.00	0.00	0.00	0.07	0.00	0.00	0.00	0.00
	0.00	0.00	0.00		0.00	0.00	0.00	0.00
	0.00	0.00	0.00		0.00	0.00	0.00	0.00
	0.00	0.00	0.00		0.00	0.00	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
1								

door locked, key dropped

0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00

0.00

0.00

0.00 0.00 0.00

0.00 0.00 0.0

0.00 0.00 0.00 0.00

0.00 0.00

0.00 0.00

0.00

beta = 0.01door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.10	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00		0.01	0.01	0.00	0.00	0.00	0.00
0.00		0.01		0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

beta = 0.001door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00		0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.17	0.00	0.00	0.00	0.00
0.00	0.00	0.00		0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

door unlocked, key carry

0.00         0.00         0.00         0.00         0.00         0.00         0.00         0.00           0.00         0.								
0.00         0.00 <td< th=""><th>0.00</th><th>0.00</th><th>0.00</th><th>0.00</th><th>0.00</th><th>0.00</th><th>0.00</th><th>0.00</th></td<>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00         0.00         0.00         0.00         0.00         0.00         0.00         0.00           0.00         0.	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00         0.00         0.00         0.06         0.06         0.06         0.11         0.00           0.00         0.	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.0	0.00	0.00	0.00					0.00
0.00 0.00 0.00 0.00 0.00 0.00 <mark>0.04</mark> 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.0	0.00	0.00	0.00	0.00	0.00	0.00		0.00
0.00 0.00 0.00 0.00 0.00 0.00 0.00	0.00	0.00	0.00	0.00	0.00	0.00		0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

door unlocked, key carry

0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.01 0.01 0.00 0.00 0.00 0.00 0.00 0.00 0.01 0.01 0.00 0.00 0.00 0.00 0.04 0.04 0.04 0.0e

0.00 0.00 0.00 0.00 0.00 0.01

0.00 0.00 0.00 0.00 0.00 0.01 0.02

0.00

0.00

0.00

PPO door locked, key carry



#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.06			0.10	0.00
0.00	0.00	0.00	0.00	0.00	0.01		0.00
0.00	0.00	0.00	0.00	0.00			0.00
0.00				0.00	0.00	0.00	



0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.00	0.00	0.00	0.00	0.00
0.00			0.02	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.04	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

beta = 0.01 door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	

beta = 0.001 door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

door locked, key dropped

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00		0.01	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.02		0.03	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

door locked, key dropped



do	PPO door locked key carry											
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00					
0.00				0.00	0.00	0.00	0.00					
0.00				0.00	0.00	0.00	0.00					
0.00			0.03	0.00	0.00	0.00	0.00					
0.00				0.00	0.00	0.00	0.00					
0.00	0.03		0.03	0.00	0.00	0.00	0.00					
0.00	0.04	0.02	0.04	0.00	0.00	0.00	0.00					
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00					

#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.01	0.00		0.00			0.00
0.00	0.00	0.00	0.00	0.00		0.01	0.00
0.00	0.00	0.01		0.00			0.00
0.00		0.01		0.02			0.00
0.00	0.00	0.00	0.01	0.00	0.01	0.00	0.00
0.00		0.00		0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00		0.01		0.00			0.00
0.00		0.01		0.00			0.00
0.00		0.01		0.00	0.01		0.00
0.00							0.00
0.00		0.01		0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.01	0.00		0.00	0.00	0.00	0.00
0.00	0.01	0.00		0.00	0.00	0.00	0.00
0.00		0.01		0.01	0.00	0.00	0.00
0.00		0.01		0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Figure 8: State distribution for standard DoorKey (1000 steps) without reward after 100k steps of training PPO

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.00	0.00	0.00	0.00	0.00
0.00			0.02	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00			0.04	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

beta = 0.01 door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	

beta = 0.001 door locked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

door locked, key dropped

0.00							
0.00	0.01			0.00	0.00	0.00	0.00
0.00		0.01	0.00	0.00	0.00	0.00	0.00
0.00	0.01			0.00	0.00	0.00	0.00
0.00	0.02			0.00	0.00	0.00	0.00
0.00	0.02			0.00	0.00	0.00	0.00
0.00	0.02			0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door locked, key dropped



PPO								
door locked, key carry								
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
0.00	0.05	0.03	0.05	0.00	0.00	0.00	0.00	
0.00	0.03			0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00			0.02	0.00	0.00	0.00	0.00	
0.00				0.00	0.00	0.00	0.00	
0.00			0.04	0.00	0.00	0.00	0.00	
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	

#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.01	0.00		0.00			0.00
0.00	0.00	0.00	0.00	0.00		0.01	0.00
0.00	0.00	0.01		0.00			0.00
0.00		0.01		0.02			0.00
0.00	0.00	0.00	0.01	0.00	0.01	0.00	0.00
0.00		0.00		0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00		0.01		0.00			0.00
0.00		0.01		0.00			0.00
0.00		0.01		0.00	0.01		0.00
0.00							0.00
0.00		0.01		0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

#### door unlocked, key carry

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.01	0.00		0.00	0.00	0.00	0.00
0.00	0.01	0.00		0.00	0.00	0.00	0.00
0.00		0.01		0.01	0.00	0.00	0.00
0.00		0.01		0.00	0.00	0.00	0.00
0.00				0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Figure 9: State distribution for standard DoorKey (1000 steps) without reward after 100k steps of training PPO



# beta = 0.01



0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

















Figure 14: State distribution for modified DoorKey (100 steps) without reward after 500k steps of training PPO

Hyperparameter	Value
PPO mini batch size	64
PPO iterations	10
PPO clipping factor ( $\epsilon$ )	0.2
PPO value coefficient $(c_1)$	0.5
PPO entropy coefficient $(c_2)$	0.003
Discount factor ( $\gamma$ )	0.99
VQ-VAE epochs	8
VQ-VAE embedding vectors $(k)$	256

Table 2: Hyperparameters for MiniGrid DoorKey

determined, the agent moves with a step size of 0.1 in the chosen direction (-0.1 for left and +0.1 for 404 405 right). Additionally, Gaussian noise with zero-mean and variance of 0.01 is added to the movement. Choosing the left action moves the agent deterministically to the left, whereas choosing the right 406 407 action moves the agent stochastically to the right, left, or keeps it in place. The left region of the 408 river is defined as  $s \in [0, 0.05]$ , where selecting the left action yields a reward of r = 5. The right region of the river is defined as  $s \in [0.95, 1]$ , where the agent receives a high reward of 10,000 if 409 410 it chooses the right action and remains in the same region. The problem is framed as a continuing 411 problem with  $\gamma = 0.99$ , where the optimal policy is to always choose the right action, while the 412 suboptimal policy results in getting stuck on the left side to collect the distractor reward. The full 413 dynamics of the environment are shown in Figure 15.



Figure 15: Riverswim environment. Agent starts at  $S_1$  and each tuple is (Action, probability, re-ward).



Figure 16: Continuous Riverswim environment. The low-reward region,  $state \in [0.00, 0.05]$ , where the agent receives r = 5 by moving left and the high-reward region, s = [0.95, 1.00] where the agent receives r = 10,000 by moving right

#### 414 S3.2 Mountain Car

415 Mountain Car (Moore, 1990) is a classic control environment consisting of a car positioned between 416 two hills. The goal is to reach the top of the right hill. To achieve this, the car must build momentum by oscillating between the hills. Mountain Car has a continuous state space, defined by the car's 417 position,  $s_1 \in [-0.12, 0.6]$ , and velocity,  $s_2 \in [-0.07, 0.07]$ . The action space is  $a \in \{$ accelerate 418 left, accelerate right, don't accelerate}. To encourage faster solutions, the agent receives a reward 419 420 of -1 at each step. The episode terminates under one of two conditions: (1) the agent reaches the 421 top of the hill  $(s_1 \ge 0.45)$ , or (2) the episode length reaches 200 steps. This setup makes Mountain 422 Car particularly challenging, as it is a sparse reward environment; therefore, the agent must explore 423 effectively to reach the top of the hill. Additionally, the short episode length provides a limited 424 window for the agent to explore efficiently.