

Exploring Non-Markov Environments Using Random Recurrent Memories

Anonymous authors

Paper under double-blind review

Abstract

1 Exploring an environment requires a reinforcement learning agent to keep track of what it
2 has already explored, which is difficult when observations do not fully reveal environment
3 state. In the Markov setting, exploration algorithms have focused on achieving systematic
4 coverage of observed states. These same methods have frequently been applied to the
5 non-Markov setting, where they aim to achieve coverage over observations. However,
6 decision-making in the non-Markov setting often depends on the agent’s entire history,
7 as opposed to only single observations. Unfortunately, achieving systematic coverage
8 over the space of all trajectories is untenable: exploration over histories is exponentially
9 expensive due to the dependence of the search space on the horizon. We therefore
10 propose a new family of methods to featurize the agent’s history with random recurrences.
11 This produces finitely-sized random statistics, or *random recurrent memories*, over an
12 agent’s history, and we aim for coverage over these memories. We describe desirable
13 properties for efficient history compression with random recurrences and propose a new
14 architecture type, the tangent recurrent unit (TRU). We show that in a diverse suite of
15 partially observable exploration tasks, tangent recurrent units, as well as other structured
16 random recurrences, outperform popular methods that aim for observation coverage.

17 1 Introduction

18 Reinforcement learning agents need to consider *what* to explore over. Many things may be informative
19 to sample for a return-maximizing agent, such as its transitions (Pathak et al., 2017; Tao et al., 2020;
20 Henaff et al., 2022) or future tests (Littman & Sutton, 2001; Ramesh et al., 2022). The most popular
21 form of exploration is to explore over the Markovian states of the environment (Bellemare et al.,
22 2016; Burda et al., 2019; Lobel et al., 2023). In the more general setting under partial observability,
23 we do not have access to the Markovian state and an agent must explore over other sampled statistics.
24 Instead of states, an agent receives partial observations at every time step, and must condition on
25 its history to inform decisions. Oftentimes Markovian exploration approaches have been applied to
26 non-Markovian observations (Burda et al., 2019; Yin et al., 2021; Zamboni et al., 2024) with limited
27 success. In this case, history information is discarded and not considered for exploration.

28 We investigate methods to explore efficiently over histories. Exploration over trajectories is difficult
29 because it requires an agent to search over *sequences* rather than only single step statistics like state.
30 This hardness is reflected when trying to take counts over trajectories; the space we are counting over
31 grows exponentially with length of the trajectory. We propose to tackle trajectory exploration with
32 structured, random fixed recurrences as a compression tool for exploration over histories. Since it can
33 always be arbitrarily difficult to compress histories in POMDPs, we must rely on properly structured
34 recurrences to efficiently disambiguate between histories. We propose parameterized recurrences
35 with properties that will allow an agent to efficiently condense down histories to explore in the
36 partially observable deep reinforcement learning setting. Based on these properties, we introduce the
37 tangent recurrent unit (TRU): a simple and principled linear recurrent function that allows agents to
38 randomly (and efficiently) encode histories over a fixed parameterization. We use these random fixed
39 recurrences as targets for density estimation to generate intrinsic rewards for deep reinforcement

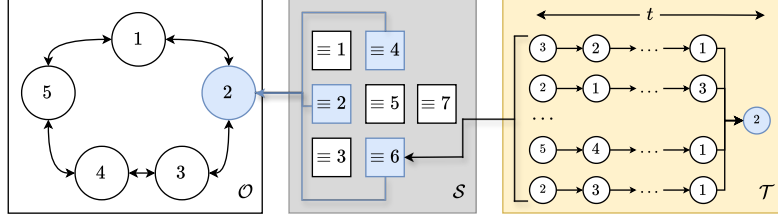


Figure 1: The Clock environment. **(Left)** Observation space exploration only gives coverage over the 5 observations. **(Center)** In state space, since observations are aliased (multiple states can produce observation 2), observation-based exploration leads to a coarse level of exploration over the modulo ground-truth states. **(Right)** Raw trajectories partition the exploration space too finely, potentially requiring an exponential (in time) number of samples for coverage.

40 learning algorithms. Across a suite of difficult partially observable exploration tasks, structured
 41 random recurrences across histories generally outperform baseline algorithms that rely on novelty
 42 over observations, or model-based novelty estimates. These experiments demonstrate the need for
 43 algorithms to explore over histories.

44 2 Background

45 We consider the problem setting of Markov decision processes (MDPs) (Puterman, 1994) and its
 46 partially observable analog (Åström, Karl Johan, 1965; Kaelbling et al., 1998). An MDP is defined
 47 by its state space \mathcal{S} , action space \mathcal{A} , reward function $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, state transition function
 48 $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta\mathcal{S}$, initial state distribution $p_0 \in \Delta\mathcal{S}$ and discount factor $\gamma \in [0, 1]$. States have
 49 the property that they are Markov: the state and action are a sufficient statistic to predict the next
 50 state and reward $P(s_{t+1}, r_t | s_t, a_t) = P(s_{t+1}, r_t | s_t, a_t, \dots, s_0, a_0)$. An agent enacts its policy
 51 $\pi_{\mathcal{S}} : \mathcal{S} \rightarrow \Delta\mathcal{A}$ to maximize its discounted future return $g_t = \sum_{i=0}^{\infty} \gamma^i r_i$, where $r_i = R(s_i, a_i)$ with
 52 $s_i \in \mathcal{S}, a_i \in \mathcal{A}$. The expectation of which is the value function $V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$. Following
 53 policy $\pi_{\mathcal{S}}$, an agent observes a distribution over states $\mu_{\pi_{\mathcal{S}}} \in \Delta\mathcal{S}$.

54 A partially observable Markov decision process (POMDP) extends MDPs by defining an observation
 55 function $\Phi : \mathcal{S} \rightarrow \Delta\mathcal{O}$. In general, observations $o \in \mathcal{O}$ are not Markov, and a statistic of the
 56 past observations and actions is required for accurate future predictions. We define sequences of
 57 observation-action pairs as the space of all possible trajectories $\mathcal{T} = (\mathcal{O} \times \mathcal{A})^* = \bigcup_{n=0}^{\infty} (\mathcal{O} \times \mathcal{A})^n$.
 58 As a shorthand, we denote an observation-action pair at time t as $x_t = (o_t, a_t)$ with $x_t \in \mathcal{X} = \mathcal{O} \times \mathcal{A}$.
 59 In POMDPs, an agent can either condition on its current observation to select actions $\pi_{\Omega} : \Omega \rightarrow \Delta\mathcal{A}$,
 60 or they can map an observed trajectory up to time t , $\tau_t \in \mathcal{T}$, to a next action $\pi : \mathcal{T} \rightarrow \Delta\mathcal{A}$ to
 61 maximize return. In order to efficiently condition over trajectories, we consider recurrences.

62 A recurrent function is a function over sequences with a recursive hidden state update. If the
 63 recurrent hidden state and input are in the sets \mathcal{H} and \mathcal{X} respectively, then a recurrent function
 64 F is defined as $F : \mathcal{H} \times \mathcal{X} \rightarrow \mathcal{H}$. Usually the recurrent hidden state space is over a fixed-size
 65 representation, like a vector over \mathbb{R}^d , to compress varying-length trajectories to a fixed size. Recurrent
 66 networks (Amari, 1972; Hochreiter & Schmidhuber, 1997; Cho et al., 2014; Jing et al., 2017) are
 67 a class of parameterized recurrent estimators that compress input sequences into hidden states that
 68 are fixed-length vectors: $f_{\theta}(x_t, h_t) = h_{t+1}$, $h_t, h_{t+1} \in \mathbb{R}^d$. These estimators have been used
 69 extensively for state estimation (Bakker, 2001; Ni et al., 2022). In this work we consider how to use
 70 these function approximators for exploration in POMDPs.

71 3 Exploration under partial observability

72 An agent must decide what to explore over and how to do it. For example, coverage over the possible
 73 transitions an agent can experience is one option for exploration, and ensures that an agent experiences
 74 transitions that may lead to an improvement in its policy. This often incurs a heavy computational cost

75 in modeling the environment (Pathak et al., 2017; Tao et al., 2020). Another option is count-based
 76 exploration (Bellemare et al., 2016; Ostrovski et al., 2017; Lobel et al., 2023): an agent is rewarded
 77 according to its epistemic uncertainty with regards to a state. This ensures coverage over the space of
 78 states.

79 Under partial observability, unless information on the belief over underlying Markovian states is
 80 available, exploration is either underspecified or requires samples that are exponential in the horizon.
 81 We make no assumptions about privileged access to information about the POMDP; we are only
 82 afforded the trajectory sequence we sample from the environment. In the Markov setting, since
 83 states are a sufficient statistic of an agent’s past, exploration over the state space is enough to afford
 84 performance in exploration problems. However, under partial observability, an agent is also required
 85 to reason about its past for better decision making, and will likely need to explore over different
 86 trajectories to improve its policy.

87 Consider the Clock problem in Figure 1. In this fixed-horizon POMDP, observations are integers
 88 from 1 to 5, the two actions (+1 and -1) stochastically transition to either adjacent observations (with
 89 a loop between 5 and 1), and an agent gets a reward for observing a total sum that is divisible by 7. In
 90 order to maximize return, an agent must try out different trajectories which lead to different modulo 7
 91 remainders. Let us consider count-based exploration in this environment by taking counts over state,
 92 observation, and trajectories.

93 **Counts over the state space \mathcal{S} .** For exploration in the Markov setting, an agent is incentivized
 94 to reach states it has not often seen. If $N(s)$ counts the number of times an agent has visited
 95 state $s \in \mathcal{S}$, then the agent sees an exploration bonus proportional to $N(s)^{-1/2}$. For count-based
 96 exploration over states in the Clock problem, an agent should try to visit sums with modulo 1 to 7, as
 97 shown in the center figure in Figure 1. In the non-Markov setting, without access to privileged state
 98 information (e.g. belief distributions), we have two alternatives: counts over observations and counts
 99 over trajectories.

100 **Counts over the observation space \mathcal{O} .** One can also consider incentivizing exploration by counting
 101 over visits to different observations. Comparing the left and center images in Figure 1, multiple
 102 states can map to the same observation. This means counting over observations is a fundamentally
 103 lossy way of counting over states—observing 2 does not guarantee the agent has visited all states
 104 that map to observation 2. In general, count-based exploration methods (even for POMDPs) aim for
 105 coverage of the state space. While using observation counts as an approximation of state counts can
 106 be beneficial in certain cases, we will show, through experimentation, that they often only work well
 107 in cases where observations themselves reveal much of the state.

108 **Counts over the space of trajectories \mathcal{T} .** While counts over trajectories is a more fine-grained
 109 measure than over observations, it can partition the input space too finely to be helpful for exploration.
 110 The right image in Figure 1 represents all the trajectories that could map to a single ground-truth state
 111 given it ends in observation 2.

112 If an agent needs to explore over all trajectories, the space that it needs to take counts over grows larger
 113 with the length of the trajectories. We can see this with visitation counts. The number of timesteps
 114 required to visit all trajectories of horizon T is bounded below by $|\mathcal{O}|^T |\mathcal{A}|^T$. With increasing
 115 trajectory lengths, this lower bound increases exponentially with the horizon T . Exploration over
 116 trajectories requires an exponentially large (in the horizon) number of samples without some form
 117 of compression over trajectories. In the next section we discuss different approaches to tackle this
 118 exponential blow up.

119 4 Compressing histories with random recurrences

120 Efficient exploration under partial observability requires a compression of trajectories to pare down
 121 the search space. Due to the nature and hardness of POMDPs, this reduction necessarily comes
 122 with trade-offs depending on the nature of the POMDP. One can construct a combinatorial lock
 123 POMDP that requires exploration over all trajectories, with no exponential blow up in exploring

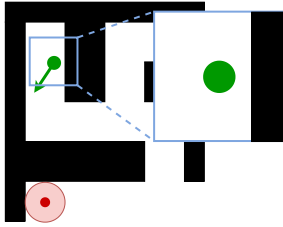
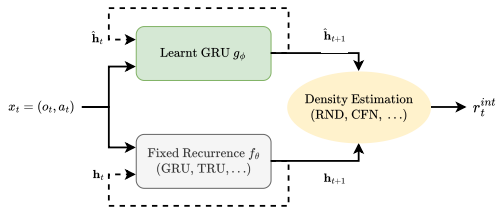


Figure 2: Visual PointMaze.

Figure 3: A general recurrent density estimation module for trajectory-based intrinsic reward r_t^{int} .

124 the state space (Liu et al., 2022). Nevertheless, there are reasonable priors we should assume about
 125 general POMDPs, and desirable properties of trajectory compressors for exploration that we now
 126 consider.

127 **Recurrences for history compression.** Instead, we can compress history so that approximate
 128 pseudocounts \hat{N} (Bellemare et al., 2016) can generalize over the compression of history, and does
 129 not grow arbitrarily large with trajectory length. One form of history compression is through
 130 recurrence. In this work, we propose to use parameterized linear and non-linear recurrent networks
 131 $h_{t+1} = f_\theta(h_t, x_t)$ to compress trajectories into fixed-length histories. These fixed-length histories
 132 allow us to define pseudocounts over trajectories.

133 **Non-stationary sequence compression in training.** Most parameterized recurrent networks are
 134 learned dynamical systems like an RNN. The issue with learned systems is that their trajectory
 135 representation (e.g. the hidden state vector in an RNN) is non-stationary over training time. Trying to
 136 calculate pseudocounts over a moving representation can lead to ill-conditioned counts. To maintain
 137 stationarity, we propose to use fixed functions of histories, i.e. fixed recurrences.

138 **Efficient disambiguation between trajectories.** A memory compressor must group together
 139 histories to avoid paying the compounding cost over time in trajectory exploration with pseudocounts.
 140 Since pseudocounts require a distance metric to be defined over trajectories, these trajectory distances
 141 should ideally reflect things that are close in state space. One approach is to learn a memory
 142 compressor by optimizing auxiliary tasks (Ramesh et al., 2022; Allen et al., 2024). This again
 143 suffers from non-stationarity. Another approach is to use *structured* random recurrences for density
 144 estimation. While it may seem that a random projection of history would not be an efficient
 145 compression of history, in the next section we show that, with the right priors and structure, we can
 146 achieve more efficient compression with random recurrences.

147 4.1 Structured random recurrences for density-based exploration

148 The biggest concern with random fixed networks is the need for efficient history disambiguation for
 149 general POMDPs. In fact, for any given finite resolution history compressor that can disambiguate
 150 sequences up to length L , the combination lock POMDP (Liu et al., 2022) can be scaled such that
 151 $|\mathcal{A}|^T > L$ and the history compressor cannot disambiguate between trajectories. While there do exist
 152 simple compressions that are able to theoretically disambiguate all histories (Eberhard et al., 2025),
 153 these compressors are generally unable to compress histories with more complex dynamics.

154 We propose structured random recurrences with density estimation for intrinsic rewards to solve
 155 general partially observed environments. Density estimation exploration approaches approximate
 156 the state visitation counts through prediction errors (Bellemare et al., 2016; Ostrovski et al., 2017;
 157 Burda et al., 2019; Lobel et al., 2023), and generates intrinsic rewards that incentivizes the agent to
 158 We visualize this proposed intrinsic reward module in Figure 3. This module has two parameterized
 159 models: a learnt recurrent function (in our case, a GRU) and a fixed, random recurrent function. We
 160 combine the outputs of these two functions for density estimation, which gives an intrinsic reward
 161 bonus. We consider what to use as a random fixed recurrence in this work. Since we cannot compress

162 trajectories from all possible POMDPs with a finite resolution history compression, we have to rely on
 163 reasonable priors to reduce the space of trajectories to efficiently disambiguate between trajectories
 164 for exploration.

165 Consider the visual PointMaze environment (Fu et al., 2020; Radji, 2025) in Figure 2. The goal of an
 166 agent is to reach the goal (the red dot). Its state and action space is continuous, and its observations
 167 are cropped, top-down, agent-centric images of the agent’s position. This requires localization based
 168 on the past history of observations to reach the goal position. We use this environment as an example
 169 for the different properties we might need for an efficient trajectory representation.

170 4.1.1 Fixed point under repeated inputs

171 Consider all the potential trajectories in visual PointMaze. A large portion of this space are trajectories
 172 where the agent is bumping into walls for extended periods of time. Repeated inputs generally include
 173 little information on the underlying state, as extended periods of the same input generally implies the
 174 state is not changing.

175 To avoid exploring these trajectories, we want our recurrent function to have a fixed point under
 176 the same input. For an input x , our recurrence is $\mathbf{h}_{t+1} = f_{\theta}(\mathbf{h}, x)$, and denote f_{θ}^k as applying the
 177 function $f_{\theta}(\cdot, x)$ repeatedly k times. Then f_{θ} has a fixed point under input x when

$$\lim_{k \rightarrow \infty} f_{\theta}^k(\mathbf{h}, x) = \mathbf{h}_x^*, \text{ where } \mathbf{h}_x^* = f_{\theta}(\mathbf{h}_x^*, x). \quad (1)$$

178 This mechanism leaves out certain memory mechanisms that may be important. For example,
 179 counting for many time steps.

180 4.1.2 Bounded updates across time

181 Many density estimators use prediction error from random projections of the input (Burda et al.,
 182 2019; Lobel et al., 2023) as a way of estimating visitation densities. These predictions assume that
 183 the space of neural network outputs is smooth in the inputs: similar state inputs should admit random
 184 representations that are close. This allows the density estimates to generalize and predict similar
 185 inputs.

186 We propose a similar mechanism for trajectory-based exploration by predicting a statistic of the
 187 random fixed recurrency. To do so, we should consider an estimator that will predict this statistic. We
 188 use a standard GRU as our estimator for density estimation. If we consider the hidden state update of
 189 a GRU for each time step, that will bound how much the hidden state of a random fixed recurrent
 190 function can change so that our predictor can match the output. Due to the final tanh activations in a
 191 GRU, we have:

$$\|\mathbf{h}' - \mathbf{h}\|_2 \leq 2\sqrt{d},$$

192 where \mathbf{h}' and \mathbf{h} are two subsequent hidden states. This implies our random fixed recurrency should
 193 produce subsequent hidden states with distances that are less than $2\sqrt{d}$.

194 We would like updates to be well below this loose bound to prevent wildly swinging predictions
 195 across time. While a GRU as the random fixed recurrent network would satisfy this bound, in the next
 196 sections, we introduce a recurrent network architecture with tighter bounds on each update, which we
 197 show can improve exploration performance.

198 We consider what forms of recurrent units have bounded updates and have a fixed point under repeated
 199 inputs. We can show that gated recurrent units (GRUs) approximate these properties. For a more
 200 principled approach, we introduce a new recurrent unit that is a simple, principled recurrency and
 201 guarantees these properties.

202 4.2 Gated recurrent unit

203 GRUs are a natural candidate for random fixed recurrent networks for exploration. As a reminder, we
 204 want to use the prediction error between the output of a frozen GRU and a learned predictor network
 205 for density estimation. In this set up, if we use a model of the same architecture and size for both

206 frozen and prediction networks, we are guaranteed our estimation model can at least *represent* the
 207 outputs of our random fixed recurrent model. This set up also trivially guarantees the fixed recurrence
 208 is bounded enough for prediction. However, we show in later sections that they are not necessarily
 209 the best for producing intrinsic rewards through density estimation.

210 Newly initialized GRUs are not guaranteed to have a fixed point under repeated inputs. Because of
 211 the orthogonal initialization of the recurrent kernel, we do not have a global contraction for fixed
 212 inputs as the spectral radius of the recurrent kernel is 1. The only guarantee we have is that a fixed
 213 point exists (Brouwer, 1911), but not necessarily will converge to it given any input.

214 4.3 Tangent recurrent unit

215 We introduce the tangent recurrent unit (TRU): a principled recurrent unit based on tangent-centered
 216 updates. TRUs are a parameterized recurrent network $\mathbf{h}_{t+1} = f_\theta(\mathbf{h}_t, \mathbf{x}_t)$ which updates its recurrent
 217 state vector (also called hidden state) with an input \mathbf{x}_t . This function happens in three steps:

$$218 \quad \mathbf{z}_t = \tilde{\mathbf{B}}\mathbf{x}_t \quad (2)$$

$$219 \quad \mathbf{z}_t^\perp = (\mathbf{I} - \mathbf{h}_{t-1}\mathbf{h}_{t-1}^\top)\mathbf{z}_t \quad (3)$$

$$220 \quad \mathbf{h}_t = \frac{\mathbf{h}_{t-1} + \mathbf{z}_t^\perp}{\|\mathbf{h}_{t-1} + \mathbf{z}_t^\perp\|_2} \quad (4)$$

218 where $\tilde{\mathbf{B}} \in \mathbb{R}^{d \times k}$ is a column-normalized. First, in Equation 2, we map the input $\mathbf{x}_t \in \mathbb{R}^k$ into
 219 \mathbb{R}^d with a parameterized linear map $\tilde{\mathbf{B}}$. Next, in Equation 3, we project the mapped input \mathbf{z}_t onto
 220 the orthogonal complement of the previous hidden state \mathbf{h}_{t-1} . Finally, in Equation 4, we add this
 221 projected vector \mathbf{z}_t^\perp to the previous hidden state and normalize to produce the new hidden state. We
 222 visualize a step of this operation in Figure 5.

223 TRUs have many beneficial properties conducive for random fixed recurrent history compression.
 224 TRUs have a bounded update not dependent on hidden dimension ($\|\mathbf{h}' - \mathbf{h}\|_2 \leq 1$), due to the
 225 normalization step. With proper initialization, TRUs have a fixed point for almost all inputs, and
 226 update the hidden state in a bounded manner due to the normalization.

227 **Theorem 1.** Let $\mathbf{x} \in \mathbb{R}^k$ be an input with $\|\mathbf{x}\|_2 = 1$, and define $\mathbf{z} := \tilde{\mathbf{B}}\mathbf{x} \in \mathbb{R}^d$. If $\|\mathbf{z}\|_2 > 0$, and if

$$228 \quad \mathbf{h}_+^* = \frac{\mathbf{z}}{\|\mathbf{z}\|_2} \quad \text{and} \quad \mathbf{h}_-^* = -\mathbf{h}_+^*,$$

228 then repeated application of f_θ reaches the fixed point \mathbf{h}_+^* for all \mathbf{x} and all initializations of $\mathbf{h}_0 \neq \mathbf{h}_-^*$:

$$229 \quad f_\theta(\mathbf{h}_+^*, \mathbf{x}) = \mathbf{h}_+^*.$$

229 *Proof.* See Section 7.1. □

230 A desirable property of recurrent functions is that the inputs are not commutative—the ordering of
 231 inputs should matter. Recurrent functions get around this with a transformation on the recurrent state;
 232 usually with a linear transformation that compounds over time. TRUs are not commutative, due to
 233 both the projection and normalization steps, since they change the hidden state with an input *relative*
 234 to the current hidden state.

235 Next, we conduct experiments to show that random fixed recurrences aid in exploration under partial
 236 observability.

237 5 Experiments

238 To test trajectory exploration capabilities, we introduce a suite of reinforcement learning environments
 239 that focus on memory-based exploration. This benchmark is entirely written in JAX (Bradbury et al.,
 240 2018), allowing for fast and scalable experimentation. We test different exploration algorithms on four
 241 environments that capture diverse forms of partial observability and differing levels of exploration
 242 difficulty.

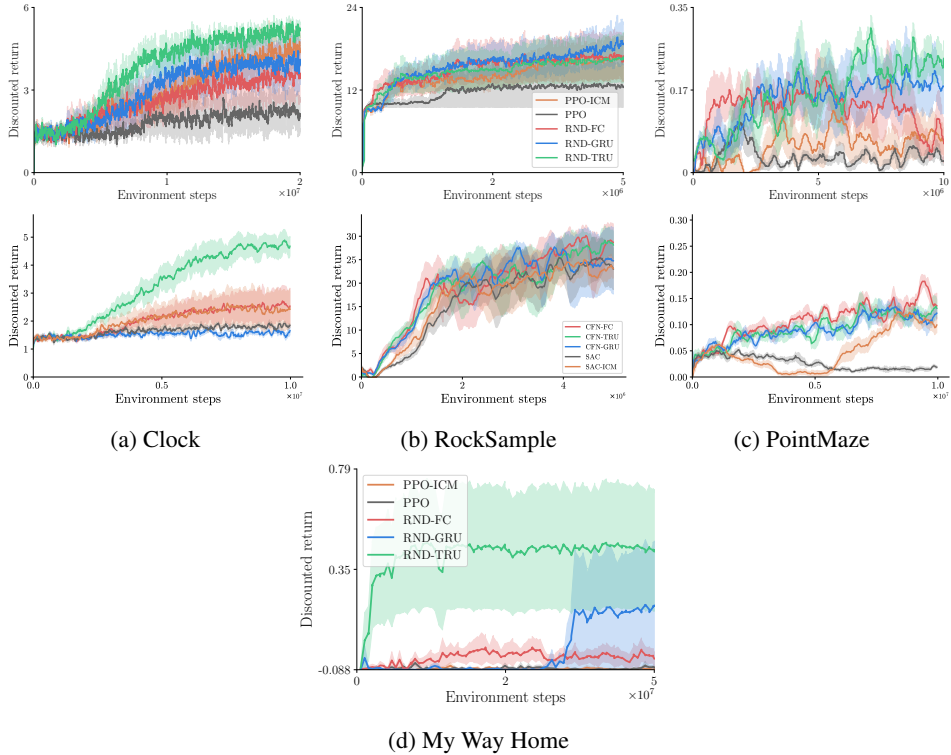


Figure 4: Performance across all four environments. Discounted episodic returns are reported for Clock, RockSample(11, 11), PointMaze, and My Way Home. The top row shows results with PPO as the base algorithm and the bottom row shows results with SAC as the base algorithm. Results are averaged across 5 seeds for all environments. Shaded regions show a 95% confidence interval.

243 **Clock** The Clock environment was introduced in Section 3. The agent receives reward at the end of
 244 an episode only when the sum of its observations is divisible by a chosen integer. Because reward
 245 depends on the full sequence rather than on any individual observation, an agent must explore over
 246 entire trajectories to discover the modular structure that yields rewards.

247 **RockSample** RockSample[n, k] (Smith & Simmons, 2004) places an agent in an $n \times n$ grid
 248 containing k rocks of unknown quality (good or bad). The agent observes its own position and
 249 receives noisy sensor readings about rock quality. It must navigate and sample the good rocks before
 250 exiting on the right edge of the grid.

251 **PointMaze** This is a continuous control navigation task (Fu et al., 2020; Radji, 2025) where an
 252 agent must traverse a maze to reach a random goal location. The observation is a 10×10 agent
 253 centric image as shown in Figure 2. Since the agent cannot localize itself from any single frame, it
 254 must do so through its history.

255 **My Way Home** A challenging visual first-person maze exploration problem introduced in ViZ-
 256 Doom (Wydmuch et al., 2019), and hardware accelerated with JAX (Bradbury et al., 2018) in
 257 JAXenstein (Tao & Konidaris, 2026). An agent randomly spawns in a location and has to navigate a
 258 maze to reach a fixed goal position.

259 5.1 Experiment setup

260 We test exploration with random fixed recurrent memories this on two exploration algorithms:
 261 proximal policy optimization (PPO) (Schulman et al., 2017) with random network distillation
 262 (RND) (Burda et al., 2019) and soft actor-critic (SAC) (Haarnoja et al., 2018) with coin flipping
 263 networks (CFN) (Lobel et al., 2023).

264 In standard instantiations of these algorithms, the exploration target is a fixed function of the current
265 observation. In our recurrent variants, we replace this target with two variations of random fixed
266 recurrences over the observation-action history: a randomly initialized GRU (RND-GRU, SAC-GRU)
267 and the tangent recurrent unit (RND-TRU, SAC-TRU) described in Sections 4.2 and 4.3 respectively.
268 The predictor network is a learned GRU trained to predict the output associated with these fixed
269 recurrent targets and we use its prediction error as an intrinsic reward.

270 We compare these against three baselines. First, we run the base algorithm without any exploration
271 bonus (PPO, SAC). Second, we run the exploration algorithms where the intrinsic reward is computed
272 with only the current observation-action pair (RND-FC, CFN-FC). Third, we compare against a
273 recurrent intrinsic curiosity module (ICM) (Pathak et al., 2017) added to PPO and SAC (PPO-ICM,
274 SAC-ICM). We sweep intrinsic reward scales and algorithm-specific loss coefficients for each method,
275 and report performance over 5 random seeds in Figure 4. Details of all experiments are given in
276 Section 8.

277 5.2 Results

278 Figure 4 shows the performance of all methods on our exploration benchmark. Overall, algorithms
279 that explore histories consistently outperform the ones that only explore over partially observed
280 observations. The only environment where the observation-based exploration matches performance
281 with the memory-based counterparts is in RockSample(11, 11). This is likely due to the structure
282 of RockSample observations: by themselves they reveal much about the state of the system. Clock,
283 PointMaze and My Way Home are three environments which all require trajectory exploration as
284 their observations do not reveal enough about state. TRUs perform particularly well in both the Clock
285 and My Way Home environments, and has matching performance on PointMaze, suggesting that it
286 may be a better candidate for trajectory-focused exploration.

287 6 Conclusion

288 Exploration over an agent’s history is necessary in most interesting cases of partial observability. We
289 present a family of methods to featurize the agent’s history with recurrences, producing fixed-sized
290 random recurrent memories for an agent to explore over. To incorporate this into a reinforcement
291 learning algorithm, we propose to use these random recurrent memories for density estimation to
292 output intrinsic rewards. For this history compression to be viable with density estimation, we test
293 recurrent structures that are both bounded and have a fixed point over repeated inputs. With these
294 properties in mind, we introduce a more principled recurrent function, the tangent recurrent unit, that
295 gives guarantees on these properties. Through deep reinforcement learning experiments, we show
296 that conditioning on trajectories in these well-known exploration algorithms leads to more effective
297 exploration in hard partially observable exploration tasks.

298 References

- 299 Cameron Allen, Aaron Kirtland, Ruo Yu Tao, Sam Lobel, Daniel Scott, Nicholas Petrocelli, Omer
300 Gottesman, Ronald Parr, Michael L. Littman, and George Konidaris. Mitigating partial observabil-
301 ity in sequential decision processes via the lambda discrepancy. In Advances in Neural Information
302 Processing Systems, volume 37, 2024.
- 303 Shun-ichi Amari. Learning patterns and pattern sequences by self-organizing nets of threshold
304 elements. IEEE Transactions on Computers, C-21(11):1197–1206, 1972.
- 305 Bram Bakker. Reinforcement learning with long short-term memory. In T. Dietterich, S. Becker,
306 and Z. Ghahramani (eds.), Advances in Neural Information Processing Systems, volume 14. MIT
307 Press, 2001.
- 308 Marc G. Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi
309 Munos. Unifying count-based exploration and intrinsic motivation. In Proceedings of the 30th
310 International Conference on Neural Information Processing Systems, NIPS’16, pp. 1479–1487,
311 2016.

- 312 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Yash Katariya, Chris Leary,
313 Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne,
314 and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL
315 <http://github.com/jax-ml/jax>.
- 316 Luitzen EJ Brouwer. Über abbildung von mannigfaltigkeiten. *Mathematische Annalen*, 71(1):97–115,
317 1911.
- 318 Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network
319 distillation. In *International Conference on Learning Representations*, 2019.
- 320 Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger
321 Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder–decoder for
322 statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in*
323 *Natural Language Processing (EMNLP)*, pp. 1724–1734, 2014.
- 324 Onno Eberhard, Michael Muehlebach, and Claire Vernade. Partially observable reinforcement
325 learning with memory traces. In *Forty-second International Conference on Machine Learning*,
326 2025.
- 327 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep
328 data-driven reinforcement learning, 2020.
- 329 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
330 maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference*
331 *on Machine Learning*, 2018.
- 332 Mikael Henaff, Roberta Raileanu, Minqi Jiang, and Tim Rocktäschel. Exploration via elliptical
333 episodic bonuses. In *Advances in Neural Information Processing Systems*. Curran Associates,
334 Inc., 2022.
- 335 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):
336 1735–1780, 1997.
- 337 Li Jing, Yichen Shen, Tena Dubcek, John Peurifoy, Scott Skirlo, Yann LeCun, Max Tegmark, and
338 Marin Soljačić. Tunable efficient unitary neural networks (eunn) and their application to rnns. In
339 *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML'17*,
340 pp. 1733–1741. JMLR, 2017.
- 341 Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra. Planning and acting in
342 partially observable stochastic domains. *Artificial Intelligence*, 101(1):99–134, 1998. ISSN
343 0004-3702.
- 344 Michael Littman and Richard S Sutton. Predictive representations of state. In T. Dietterich, S. Becker,
345 and Z. Ghahramani (eds.), *Advances in Neural Information Processing Systems*, volume 14. MIT
346 Press, 2001.
- 347 Qinghua Liu, Alan Chung, Csaba Szepesvari, and Chi Jin. When is partially observable reinforcement
348 learning not scary? In Po-Ling Loh and Maxim Raginsky (eds.), *Proceedings of Thirty Fifth*
349 *Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pp.
350 5175–5220. PMLR, 02–05 Jul 2022.
- 351 Sam Lobel, Akhil Bagaria, and George Konidaris. Flipping coins to estimate pseudocounts for
352 exploration in reinforcement learning. In *Proceedings of the 40th International Conference on*
353 *Machine Learning, ICML'23*, 2023.
- 354 Chris Lu, Jakub Kuba, Alistair Letcher, Luke Metz, Christian Schroeder de Witt, and Jakob Foerster.
355 Discovered policy optimisation. *Advances in Neural Information Processing Systems*, 35:16455–
356 16468, 2022.

- 357 Tianwei Ni, Benjamin Eysenbach, and Ruslan Salakhutdinov. Recurrent model-free RL can be
358 a strong baseline for many POMDPs. In International Conference on Machine Learning, pp.
359 16691–16723. PMLR, 2022.
- 360 Georg Ostrovski, Marc G. Bellemare, Aäron van den Oord, and Rémi Munos. Count-based exploration
361 with neural density models. In Proceedings of the 34th International Conference on Machine
362 Learning - Volume 70, ICML’17, pp. 2721–2730. JMLR, 2017.
- 363 Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by
364 self-supervised prediction. In International Conference on Machine Learning (ICML), 2017.
- 365 Martin L. Puterman. Markov Decision Processes. Wiley, 1994.
- 366 Waris Radji. Pointax: Jax-native pointmaze environment, 2025. URL [https://github.com/
367 riiswa/pointax](https://github.com/riiswa/pointax).
- 368 Aditya Ramesh, Louis Kirsch, Sjoerd van Steenkiste, and Jürgen Schmidhuber. Exploring through
369 random curiosity with general value functions. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave,
370 K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems, volume 35, pp.
371 18733–18748. Curran Associates, Inc., 2022.
- 372 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
373 optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.
- 374 Trey Smith and Reid Simmons. Heuristic search value iteration for pomdps. In Proceedings of the
375 20th Conference on Uncertainty in Artificial Intelligence, UAI’04, 2004.
- 376 Ruo Yu Tao and George Konidaris. Jaxenstein: Accelerated benchmarking for first-person environ-
377 ments, 2026.
- 378 Ruo Yu Tao, Vincent Francois-Lavet, and Joelle Pineau. Novelty search in representational space for
379 sample efficient exploration. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin
380 (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 8114–8126. Curran
381 Associates, Inc., 2020.
- 382 Marek Wydmuch, Michał Kempka, and Wojciech Jaśkowski. ViZDoom Competitions: Playing Doom
383 from Pixels. IEEE Transactions on Games, 11(3):248–259, 2019. The 2022 IEEE Transactions on
384 Games Outstanding Paper Award.
- 385 Haiyan Yin, Jianda Chen, Sinno Jialin Pan, and Sebastian Tschiatschek. Sequential generative
386 exploration model for partially observable reinforcement learning. Proceedings of the AAAI
387 Conference on Artificial Intelligence, 35(12):10700–10708, May 2021.
- 388 Riccardo Zamboni, Duilio Cirino, Marcello Restelli, and Mirco Mutti. The limits of pure exploration
389 in pomdps: When the observation entropy is enough. Reinforcement Learning Journal, 2:676–692,
390 2024.
- 391 Åström, Karl Johan. Optimal Control of Markov Processes with Incomplete State Information I. 10:
392 174–205, 1965. ISSN 0022-247X.