A MODEL OF PLACE FIELD REORGANIZATION DURING REWARD MAXIMIZATION

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

030

Paper under double-blind review

ABSTRACT

When rodents learn to navigate in a novel environment, a high density of place fields emerges at reward locations, fields elongate against the trajectory, and individual fields change spatial selectivity while demonstrating stable behavior. Why place fields demonstrate these characteristic phenomena during learning remains elusive. We develop a normative framework using a reward maximization objective, whereby the temporal difference (TD) error drives place field reorganization to improve policy learning. Place fields are modeled using Gaussian radial basis functions to represent states in an environment, and directly synapse to an actor-critic for policy learning. Each field's amplitude, center, and width, as well as downstream weights, are updated online at each time step to maximize cumulative reward. We demonstrate that this framework unifies the three disparate phenomena observed in navigation experiments. Furthermore, we show that these place field phenomena improve policy convergence when learning to navigate to a single target and relearning multiple new targets. To conclude, we develop a normative model that recapitulates several aspects of hippocampal place field learning dynamics and unifies mechanisms to offer testable predictions for future experiments.

1 INTRODUCTION

A place field is canonically described as a localized region in an environment where the firing rate of 032 a hippocampal neuron is maximal and robust across trials (O'Keefe, 1978; O'Keefe & Dostrovsky, 033 1971). Classically, each neuron has a unique spatial receptive field such that the population activity 034 can describe an animal's allocentric position within the environment (Moser et al., 2015). Ablation 035 studies demonstrate that the hippocampal representation is useful for learning to navigate to new targets (Morris et al., 1982; Packard & McGaugh, 1996; Steele & Morris, 1999). Importantly, each 037 field's spatial selectivity evolves with experience in a new environment before stabilizing in the later stages of learning (Frank et al., 2004). Specifically, a high density of place fields emerge at reward locations (Gauthier & Tank, 2018; Lee et al., 2020; Sosa et al., 2023), place fields elongate 040 backward against the trajectory (Mehta et al., 1997; Priestley et al., 2022), and individual field's 041 spatial selectivity continues to change or "drift" even when animals demonstrate stable behavior (Geva et al., 2023; Kentros et al., 2004; Krishnan & Sheffield, 2023; Mankin et al., 2012; Ziv et al., 042 2013). Although disparate mechanisms have been proposed to model these phenomena, a framework 043 that can unify their phenomena and clarify their computational role remains elusive. 044

Here, we propose a normative model for spatial representation learning in hippocampal CA1, given its role in representing salient spatial information (Dong et al., 2021; Dupret et al., 2010). Our primary contributions are as follows:

048

We develop a two-layered reinforcement learning model to study spatial representation learning by place fields (Fig.1A). The first layer contains a population of Gaussian radial basis functions that transform continuous spatial information into a relevant representational substrate, which feed into the actor-critic network in the second layer that uses these representations to maximize cumulative discounted reward. Besides the actor and critic weights, each place field's firing rate, center of mass and width is optimized by the temporal difference error.

- Our model recapitulates three experimentally-observed neural phenomena during task learning: the emergence of high place field density at rewards, elongation of fields against the trajectory, and drifting fields that do not affect task performance.
- We analyze the factors that influence these representational changes: a low number of fields drives greater spatial representation learning, each place field's firing rate reflects the value of that location, and increasing noise magnitude during field parameter updates causes a monotonic decrease in population vector correlation but non-monotonic change in behavior.
 - We demonstrate that optimizing place field widths and amplitudes enhances reward maximization and policy convergence. However, field parameter optimization alone is insufficient for learning to navigate to new targets. Introducing noisy field parameter updates improves new target learning, suggesting a functional role for noise.
- 064 065 066 067

061

062

063

2 RELATED WORKS

069 Anatomically constrained architecture for navigation. Learning to navigate involves the hippocampus encoding spatial information and its strong glutamatergic connections to the striatum (Floresco et al., 2001; Lisman & Grace, 2005). The ventral and dorsal regions of the striatum are 071 associated with value estimation and stimulus-response associations, functioning similarly to a critic 072 and an actor, respectively (Houk et al., 1994; Joel et al., 2002; Niv, 2009). Additionally, dopamine 073 neurons in the Ventral Tegmental Area influence plasticity in the striatal synapses (Reynolds et al., 074 2001; Russo & Nestler, 2013). This anatomical insight has led to the design of a biologically plausi-075 ble navigation model, where place fields connect directly to an actor-critic framework, and synapses 076 are modulated by the TD error (Arleo & Gerstner, 2000; Brown & Sharp, 1995; Foster et al., 2000; 077 Frémaux et al., 2013; Kumar et al., 2022). Furthermore, recent evidence shows direct dopaminergic 078 projections to the hippocampus to modulate place cell activity, strengthening the case for navigation 079 models with adaptive place fields (Kempadoo et al., 2016; Krishnan et al., 2022; Palacios-Filardo & 080 Mellor, 2019; Sayegh et al., 2024). How upstream information from the entorhinal cortex influences 081 place field representations for policy learning needs clarity (Bush et al., 2015; Fiete et al., 2008).

Field density increases near reward locations. As animals learn to navigate in a 1D track, a high density of place fields emerge at reward locations. We define density to be both the number of fields (Gauthier & Tank, 2018; Sosa et al., 2023) and the peak firing rate of each field (Lee et al., 2020).
Reward location based reorganization was observed in hippocampal CA1 and not in CA3 (Dupret et al., 2010).

087 Fields learn to encode future occupancy. As animals traverse a 1D track towards a reward, most 088 CA1 fields increase in size and their center of mass shift backwards against the trajectory of motion 089 (Frank et al., 2004; Mehta et al., 1997; Priestley et al., 2022). A proposal for this behavior is that fields initially encoding only location x_t are learning to also encode the previous location x_{t-1} , and 091 hence are coding future location occupancy $p(x_{t+1}|x_t)$ (Mehta et al., 2000; Stachenfeld et al., 2017). 092 While algorithms such as the successor representation (Dayan, 1993) learn to predict the transition 093 structure (Gardner et al., 2018; Gershman, 2018), the representation is dependent on a predefined navigation policy. Hence, a complete normative argument-including policy learning-for why 094 fields exhibit this behavior is still lacking. 095

096 Fields drift during stable behavior. After animals reach a certain performance criterion in navigat-097 ing to a reward location, the spatial selectivity of individual place fields changes across days, even 098 though animals exhibit stable behavior (de Snoo et al., 2023; Geva et al., 2023; Kentros et al., 2004; 099 Mankin et al., 2012; Ziv et al., 2013). A proposal is that these fields continue to drift within a degenerate solution space while the overall representational manifold or the chosen performance metric 100 remains stable (Kappel et al., 2015; Masset et al., 2022; Pashakhanloo & Koulakov, 2023; Qin et al., 101 2023; Rokni et al., 2007). However, a model that demonstrates stable navigation learning behavior 102 with drifting fields is absent. Furthermore, why drifting fields might be useful is still unexplored. 103

Place fields versus place cells. Several experiments have shown that place fields along the dorso-ventral axis have different widths (Jung et al., 1994) and are also involved in navigation (Contreras et al., 2018; Harland et al., 2021), while newer experiments challenge the canonical definition that a place cell only has one place field (Eliav et al., 2021). As a simple starting point, in this work we study spatial representational learning using Gaussian place fields, instead of place cells.

¹⁰⁸ 3 TASK AND MODEL SETUP

Most navigational experiments involve an animal moving from a start location to a target location to 110 receive a reward, either in a one-dimensional (1D) track or a two-dimensional (2D) arena. Similarly, 111 our agents receive their true position at every time step (t) described by the variable (scalar x_t in 1D, 112 vector x_t in 2D), and have to learn a policy (π) that specifies the actions to take (g_t) to move from 113 a start location (e.g. $x_{start} = -0.75$, Fig. 1A green dash) to a target with reward values following 114 a Gaussian distribution ($x_r = 0.5, \sigma_r = 0.05$, Fig. 1A red area). The agent outputs a discrete one-hot vector g_t (left or right in 1D and left, right, up or down in 2D), which causes its motion 115 116 to be discrete, similar to a trajectory in a grid world. To model smooth trajectories in a continuous space as an animal's behavior, we use a low-pass filter to smooth g_t using a constant $\alpha_{env} = 0.2$ 117 after scaling for maximum displacement using $v_{max} = 0.1$: 118

$$x_{t+1} = x_t + a_t$$
, $a_{t+1} = (1 - \alpha_{env})a_t + \alpha_{env}v_{max}g_t$,

similar to past works (Foster et al., 2000; Frémaux et al., 2013; Kumar et al., 2022; 2024b; Zannone et al., 2018). To track an agent's reward maximization performance during navigational learning we compute the true cumulative discounted reward ($G = \sum_{t=0}^{T} \gamma^t r_{t+1}$) for each trial using $\gamma = 0.9$ as the discount factor, which is similar to tracking the cumulative reward. The trial is terminated if the trial time reaches a threshold T_{max} or when the total reward achieved $\sum_{t=0}^{T} r_{t+1}$ reaches a threshold R_{max} . For further details, see App. A.

125 126 127

128

131

141

142

147

119 120

121

122 123 124

3.1 PLACE FIELDS AS SPATIAL FEATURES

The agent represents space through N place fields, which have spatial selectivity modeled as simple Gaussian bumps and tile the environment:

$$\phi_i(x_t) = \alpha_i^2 \exp(-||x_t - \lambda_i||_2^2 / 2\sigma_i^2),$$
(2)

(1)

where α , λ and σ set the amplitude, center, and width respectively. Two types of place field distri-132 butions were initialized: (1) Homogeneous population of constant values for amplitudes $\alpha_i = 0.5$, 133 widths $\sigma_i = 0.1$, and centers uniformly tiling the environment $\lambda = [-1, ..., 1]$ (Foster et al., 2000; 134 Frémaux et al., 2013; Kumar et al., 2022; 2024b; Zannone et al., 2018). (2) Heterogeneous popula-135 tion with amplitudes, widths and centers drawn from uniform random distributions between [0, 1], 136 $[10^{-5}, 0.1]$, [-1, 1] respectively. These ranges are consistent with experimental data where place 137 fields were 20 cm to 50 cm wide in a small environment (Frank et al., 2004; Lee et al., 2020; Mehta 138 et al., 1997; Sosa et al., 2023). 2D place fields have scalar amplitudes, two dimensional vectors for 139 center, and square covariance matrices for the width (Menache et al., 2005). Refer to App. A. 140

3.2 POLICY LEARNING USING AN ACTOR-CRITIC

To model an animal's trial-and-error based learning behavior, we adopt the reinforcement learning framework, specifically the actor-critic (Arleo & Gerstner, 2000; Brown & Sharp, 1995; Foster et al., 2000; Frémaux et al., 2013; Kumar et al., 2022; 2024b). The critic linearly weights place field activity using a vector w_i^v to estimate the value of the current location

v

$$u(x_t) = \sum_i^N w_i^v \phi_i(x_t) \,. \tag{3}$$

The value of a location corresponds to the expected cumulative discounted reward for that location. The actor has M units, each specifying a movement direction. In the 1D and 2D environments, M = 2 and M = 4 respectively to code for opposing directions in each dimension e.g. left versus right and up versus down. Each actor unit a_j linearly weights the place field activity such that the matrix W_{ji}^{π} computes the preference for moving in the j-th direction

$$a_j(x_t) = \sum_i^N W_{ji}^{\pi} \phi_i(x_t) \quad , \quad P_j = \frac{\exp(a_j)}{\sum_k^M \exp(a_k)} \,,$$
 (4)

with the probability of taking an action computed using a softmax. A one-hot vector g_j is sampled from the action probability distribution P as in Foster et al. (2000), making this policy stochastic. w_i^v and W_{ji}^{π} were initialized by sampling from a random normal distribution $\mathcal{N}(0, 10^{-5})$.

158 159

153 154

3.3 BIOLOGICALLY RELEVANT REWARD MAXIMIZATION LEARNING OBJECTIVE

The objective of our agent is to maximize the expected cumulative discounted reward $\mathcal{J}^G = \mathbb{E}[G_t] = \mathbb{E}[\sum_{k=0}^T \gamma^k r_{t+1+k}]$. To achieve this goal in an online manner, our agent uses the stan-

dard actor-critic algorithm using the expected temporal difference objective (refer to App. A):

$$\mathcal{J}^{TD} = \mathbb{E}\left[\sum_{t}^{T} r_{t+1} + \gamma v(x_{t+1}) - v(x_t)\right] = \mathbb{E}\left[\sum_{t}^{T} \delta_t\right].$$
(5)

which reduces variance and speeds up policy convergence (Dayan & Abbott, 2005; Mnih et al., 2016; Schulman et al., 2017; Sutton & Barto, 2018; Wang et al., 2018). The TD error is also biologically relevant, as the responses of midbrain dopamine neurons resemble TD reward prediction error (Amo et al., 2022; Gershman & Uchida, 2019; Montague et al., 1996; Schultz et al., 1997; Starkweather & Uchida, 2021). The actor learns a reward maximizing policy by ascending the gradient of the policy log likelihood, modulated by the TD error. To accurately estimate the TD error and critique policy learning, the critic learns a value function by minimizing the squared TD error $\mathcal{L} = \mathbb{E}\left[\sum_{t}^{T} \frac{1}{2}\delta_{t}^{2}\right]$.

As our agent uses a single population of place fields, these fields must learn spatial features that enhance both policy and value learning. The field parameters $\theta = \{\alpha, \lambda, \sigma\}$ and the policy weights W^{π}, w^{v} are updated by gradient ascent using a joint objective modified from Wang et al. (2018):

$$\nabla_{\theta, W^{\pi}, w^{v}} \mathcal{J} = \nabla_{\theta, W^{\pi}} \mathcal{J}^{TD} - \nabla_{\theta, w^{v}} \mathcal{L} = \mathbb{E} \left[\sum_{t}^{T} \left(\nabla_{\theta, W^{\pi}} \log \pi(g_{t} | x_{t}) + \nabla_{\theta, w^{v}} v(x_{t}) \right) \cdot \delta_{t} \right],$$
(6)

177 with $\nabla_{w^v} \mathcal{J}^{TD} = 0$ and $\nabla_{W^{\pi}} \mathcal{L} = 0$. We estimate all parameter gradients online, and provide the 178 explicit update equations for each parameter in App. A. The learning rates for the actor-critic and 179 place field parameters can be the same (Sup. Fig. 13). For theoretical analysis, we assume a sep-180 aration of timescales between learning the actor-critic weights and updating place field parameters 181 (App. B). This approach stabilizes place field representation learning, and is consistent with Dong 182 et al. (2021)'s observation that rodent behavior converges faster than place field representations.

183 184

185

186

176

4 Results

4.1 A HIGH DENSITY OF FIELDS EMERGES NEAR THE REWARD LOCATION

187 We first examine the neural phenomenon where a high density of place fields emerges at the reward 188 location. Field density is defined by the distribution of field centers of mass (COM) (Gauthier & 189 Tank, 2018), which we estimate using Gaussian kernel smoothing. Figure 1B shows how our agent's 190 track occupancy, field density, mean firing rate, and individual field's spatial selectivity change when 191 learning to navigate in a 1D track from the start location $x_{start} = -0.75$ to the target at $x_r = 0.5$, when only optimizing place field centers ($\Delta\lambda$). In the early stages of learning, the agent spends a 192 higher proportion of time at the start location with sporadic exploration towards the reward. Despite 193 this behavior, a high field density and mean firing rate rapidly emerges at the target from a homoge-194 neous field population within the first few trials. Individual fields at the reward location shift closer 195 to the target (Fig. 1F), as seen in (Gauthier & Tank, 2018; Sosa et al., 2023), in contrast to fields at 196 non-rewarded locations. As learning progresses and the agent spends a higher proportion of time at 197 the reward location, field density and mean firing rate at the start location also begins to rise slightly, 198 although it remains lower than at the reward location, replicating the two-peaked field distribution 199 in (Gauthier & Tank, 2018). A high density at the reward location followed by the start location ro-200 bustly emerges in heterogeneous place field populations when all the field parameters $(\Delta\lambda, \Delta\alpha, \Delta\sigma)$ 201 are optimized (Fig. 1B right, Sup. Fig. 2B). Similar field dynamics are observed in a 2D arena with an obstacle where agents have to navigate to a target from a starting location (Fig. 1C). When op-202 timizing all the field parameters in a homogeneous population, a high field density rapidly emerges 203 at the reward location to increase goal representation as seen in (Dupret et al., 2010), followed by 204 gradual reorganization of field density along the agent's trajectory back to the start location. 205

206 Interestingly, increasing the number of fields in a heterogeneous place field population reduced the average density (Fig. 1D, Sup. Fig. 1) and mean firing rate (Sup. Fig. 4D) that emerges near 207 the reward location (R > 0.01). This is because as the number of fields increase, the agent goes 208 into a weak feature learning regime (Sup. Fig. 4) in which feature learning does not contribute to 209 additional advantage. While experiments can record thousands of place fields, only a small fraction 210 of fields, between 80 to 150, show reward-relative reorganization (Gauthier & Tank, 2018; Lee et al., 211 2020; Sosa et al., 2023). Conversely, the density and mean firing rate are proportional to the reward 212 magnitude (blue versus green), and inversely proportional to the reward location width (red versus 213 purple) as a narrower target might require higher discriminability for the agent to maximize rewards. 214

To understand why place fields exhibit these dynamics, we perform a perturbative approximation to the place field parameter changes under TD learning updates (Bordelon et al., 2024; Menache et al.,



Figure 1: Emergence of high field density at the reward location with learning. (A) The task is 236 to navigate from the start (green dash) to the target (red area) to receive rewards whose magnitude 237 follows a Gaussian distribution. The agent contains N Gaussian place fields (blue) which synapse to 238 an actor (red) and critic (green) to learn the policy and value function respectively. The temporal dif-239 ference error δ is used to update parameters. (B-C) Example of an increase in place field density at 240 the reward location during learning in a (B) 1D track (Gauthier & Tank, 2018; Lee et al., 2020), and 241 (C) 2D arena (Dupret et al., 2010) with an obstacle (gray). (B) When optimizing field centers (Top 242 row) In the early learning phase (T = 100), the agent spends a high proportion of time ($p_{BM}(x)$, 243 black) at the start location with a constant field density (gKDE(COM), blue) throughout the track. 244 As learning proceeds (T = 2000, 50000), a higher field density emerges at the reward and start lo-245 cation when only optimizing field centers ($\Delta\lambda$). (Bottom row) Evolution of individual field centers and mean firing rates ($\sum \phi(x)$, red). (Right) A high field density and mean firing rate emerges at 246 the reward location, followed by the start location, for a heterogeneous place field population when 247 all field parameters are optimized ($\Delta\lambda, \Delta\alpha, \Delta\sigma$). (C) The density similarly evolves in the 2D arena 248 when all field parameters are optimized. In the early learning phase (T < 10000), centers of mass 249 (COM, black dots) shift to the target, causing a high density to emerge at the reward (right). In the 250 later learning phase, the rest of the COM align along the trajectory. The start and reward locations 251 and radius for goal representation (G.R.) are marked by green, red and blue circles in the leftmost plot. (D) As the number of fields increases, the average field density (d(x) = qKDE(COM)) near 253 the reward location x_r compared to non-reward location x' decreases for the heterogeneous popu-254 lation. The density decreases when the reward magnitude decreases ($R_{max} = 1, 5, 9$: blue, orange, 255 green) and reward location's size increases ($R_{size} = 0.025, 0.05, 0.1$: red, orange, purple). (E-F) Example of field dynamics when an agent (N = 512) navigates a 1D track. (E) Fields initialized 256 before ($\lambda_i = 0.5$, blue) and after ($\lambda_i = 0.6$, orange) the target move forward and backward respec-257 tively, increasing the density near the target. (F) Fields closest to the reward ($\lambda_i = 0.5, 0.6$: green 258 and red) show a rapid and high amplification compared to other fields ($\lambda_i = -0.75, 0.0$: blue and 259 orange). The first order perturbative prediction (theory) provides a good approximation of learned 260 amplitudes in both 1D and 2D tasks. Shaded area and error bars are 95% CI over 50 seeds. 261

264

2005). In this approximation, we assume that the change to the field parameters is small, controlled by the number of fields, and by the large separation between learning rates. Focusing on the place field centers, we derive in App. B the approximation where $\eta_{\lambda} = 0.0001$ is the learning rate for the field centers and $\eta = 0.01$ is the learning rate for the critic weights:

$$\lambda_i(t) - \lambda_i(0) \approx \frac{\eta_\lambda}{\eta} \left(\frac{2}{\sigma_i^2} + \frac{1}{\sigma_x^2}\right)^{-1} \left[\frac{\lambda - \lambda_i(0)}{\sigma_i^2} + \frac{\bar{\mu}_x - \lambda_i(0)}{\sigma_x^2}\right] w_{v,i}^2(t) , \ \eta_\lambda \ll \eta , \tag{7}$$

⁶⁹ Under this approximation, each field's center shifts proportionally to the squared magnitude of the critic weights (w_n^2) , implying that fields at locations with a high value will shift at a faster rate



Figure 2: Reward maximization predicts field enlargement against movement direction, with 283 field dynamics distinct from the successor representation. (A-B) Both Reward Maximization 284 (RM, orange) and Successor Representation (SR, blue) algorithms cause (A) field sizes to increase and (B) field center of mass to shift backwards against the movement direction when learning in a 1D track, replicating Mehta et al. (1997). Each line shows the average change in an agent initialized 287 with 16 place fields. The change in SR and RM fields were normalized separately to be between 288 0 to 1 for visualization. (C) In the early learning phase (T = 1000), both the SR (top row) and 289 RM (bottom row) agents spend a high proportion of time at the start location (black), and learn 290 a policy to spend a higher proportion of time at the target in later phases (T = 10000, 50000). 291 The individual SR fields (colored) and SR mean firing rate (red) closely track the proportion of 292 time the agent spends in a location. Conversely, the individual RM fields and mean firing rate 293 show an inverse relationship against the proportion of time the RM agent spends at a location in the early learning phase, but start to align in the later phases. (D) The proportion of time SR and RM agents spend at a location is high, positively correlated (black). SR agents show a consistently 295 high, positive correlation (blue) between mean firing rate $(\sum \psi(x))$ and proportion of time spent in 296 a location $(p_{SR}(x))$. Conversely, the correlation between the RM agents' mean firing rate $(\sum \phi(x))$ 297 and time spent at a location $(p_{RM}(x))$ becomes anti-correlated (orange) before becoming positively 298 correlated. Similarly, the SR and RM field densities (red) become anti-correlated before becoming 299 positively correlated at the later learning phase. (E) The correlation between the individual field 300 firing rates ($\psi_i(x)$ vs $\phi_i(x)$, green) and the spatial representation similarity matrices ($\psi(x) \cdot \psi(x')$) 301 vs $\phi(x) \cdot \phi(x')$, purple) learned by the SR and RM agents rapidly diverge in the early learning phase 302 but stabilize and become positively correlated in later phases. (F) Example change in field size and 303 COM by SR (top row) and RM (bottom row) agents in a 2D arena with an obstacle. Summary 304 statistics in Sup. Fig. 6. The RM agent's field elongation is more pronounced than the SR agent, especially along the trajectory and rotation about the obstacle. Shaded area is 95% CI over 10 seeds. 305

compared to locations with a low value. In addition to the value of a location, the agent's start 307 location (modeled as a Gaussian with mean $\bar{\mu}_x = -0.75$ and spread σ_x) and the mean field center 308 location λ over time under the policy influence each field's displacement. As the reward location is 309 visited frequently, we expect $\lambda \approx 0.5$. As the term within the square bracket changes sign depending 310 on the field location, only the fields near the reward location will shift towards the reward, while the 311 rest of the fields will move towards the start location. Due to these influences, the field density at 312 the reward location will increase first followed by a gradual increase in start location (Fig. 1B,E). 313 Additional approximations are needed to model the agent's trajectory and improve the simulation-314 theory fit for place field centers (App. B). A similar perturbative analysis for amplitudes yields $\alpha_i(t) - \alpha_i(0) \approx 2 \frac{\eta_\alpha}{\eta} w_{v,i}^2(t)$ when $\eta_\alpha \ll \eta$, where $\eta_\alpha = 0.0001$ is the learning rate for the α 315 316 parameters. Thus, fields at locations with a high value will be amplified at a rate similar to the agent learning the value function (Fig. 1F). Therefore, this approximation predicts fields shifting to the 317 start and reward location with field amplification at the reward location. 318

319

321

306

4.2 REWARD MAXIMIZATION RESULTS IN FIELD ENLARGEMENT AGAINST MOVEMENT

We now turn to the next phenomenon where place field sizes increase and their centers shift backward against the movement direction as animals learn to navigate. A proposed account for this phenomenon is that place fields learn to encode future occupancy, that is, given a location x_t , the population of place fields represents the future occupancy probability $p(x_{t+1}|x_t)$ (Stachenfeld et al., 2017). Future occupancy can be learned through Hebbian association of fields that have a correlated firing activity sequence (George et al., 2023; Mehta et al., 2000), or through the successor representation (SR) algorithm, whose objective is to minimize state prediction error by computing a temporal difference error for each place field to learn the transition probabilities (Dayan, 1993; Gardner et al., 2018). Both methods recapitulate field elongation in a 1D track. Here, we show that our reward maximizing (RM) agent does as well.

331 For comparison purposes, we developed an SR agent that learns the transition probabilities in par-332 allel to policy learning (Sup. Fig. 5A). The SR agent has a similar architecture to our (RM) agent 333 (Fig. 1A), with two key differences: 1) It has one set of place fields with fixed parameters, and only 334 the synapses from these place fields to the actor-critic are optimized for policy learning. 2) There is a separate set of N successor place fields $\psi(x)$ that receive input from the fixed place fields via 335 synapses U which are optimized using the SR algorithm (App. C). We will compare the learned suc-336 cessor place fields to the learned place fields in our RM model, following Stachenfeld et al. (2017). 337 We will therefore henceforth refer to the successor place fields simply as place fields. 338

339 Both SR and RM agents recapitulate the phenomena seen in (Mehta et al., 1997; Priestley et al., 340 2022): on average, place fields increase in size over learning (Fig. 2A), and the center of mass 341 (COM) shifts backwards from their initialized positions (Fig. 2B, Sup. Fig. 5C). However, the place fields of the SR and RM agents evolve differently. Both the SR and RM agents initially spend a high 342 proportion of time at the start location and gradually learn a policy to spend a higher proportion of 343 time at the reward location (Fig. 2C). The correlation between the SR and RM agents proportion of 344 time spent in a location is high, positively correlated in most trials (Fig. 2D), except for the decrease 345 between trial 5000 to 10,000 where the RM agent spends a higher proportion of time at the reward 346 location than the SR agent due to faster policy convergence (Sup. Fig. 5B). 347

The SR, by design, tracks the transition probabilities of the agent's policy. Consequently, the SR 348 mean firing rate $\sum \psi(x)$ closely aligns with the agent's probability of spending time at a location 349 p_{SR} , showing a high positive correlation (Fig. 2C, D). Conversely, during early learning, the RM 350 agent exhibits a high mean firing rate $\sum \phi(x)$ at the reward location, which contrasts sharply with 351 the time proportion spent at that location (Fig. 2C), leading to a highly negative correlation between 352 $\sum \phi(x)$ and p_{RM} (Fig. 2D). Interestingly, in the later phase of learning, $\sum \phi(x)$ and p_{RM} become 353 positively correlated. The mean firing rates learned by the SR and RM agents become negatively 354 correlated during the early learning phase but become positively correlated at the later learning phase 355 (Fig. 2D). A similar change in correlation is observed when comparing the individual SR and RM 356 field selectivity or population vectors (Fig. 2E), and the spatial representation similarity matrix (Sup. 357 Fig. 5D) by taking the dot product of SR and RM field firing rates at all locations (Fig. 2E). This 358 demonstrates that both algorithms eventually learn similar spatial representations, but the process of learning these representations are different. 359

In a 2D arena with an obstacle, both agents show elongation of fields against the agent's direction of movement (Fig. 2, Sup. Fig. 6) while also accounting for the blockage of path by the obstacle. The RM agent shows a significantly larger elongation of fields to span the entire corridor while the elongation of fields by SR is subtle.

364 365

366

4.3 STABLE NAVIGATION BEHAVIOR WITH DRIFTING FIELDS

367 The third phenomena that the model captures has been described as representational drift, where 368 the agent demonstrates stable behavior but the spatial selectivity of individual place fields changes 369 over time (Fig. 3A, Sup. Fig. 8G). Although our agent uses a stochastic policy, both the navigation 370 behavior after 25,000 trials (Fig. 3C, blue) and the population vector (PV) correlation are extremely 371 stable (Fig. 3B, blue). To drive larger variability in the representation, we introduced Gaussian noise 372 to the field parameter updates at every time step (App. D). Increasing the noise magnitude led to 373 a faster decrease in PV correlation but also disrupted agents' policy convergence for magnitudes 374 greater than 10^{-3} (Sup. Fig. 7). Hence, we consider the noise magnitudes between 10^{-4} and 10^{-3} . As the noise magnitude increases, agent's reward maximization behavior remains stable 375 376 while the PV correlation decreases rapidly (Fig. 3B-C). This demonstrates that agents can optimize their policies to maintain stable behavior even though individual spatial selectivity is changing. 377 Interestingly, the spatial representation similarity matrix remains more stable than PV correlation



Figure 3: Stable behavior and representation similarity despite drifting fields. (A) Inject-397 ing Gaussian noise with magnitude $\sigma_{noise} = 0.0001$ into field parameters causes individual 398 field's spatial selectivity to change across trials. (B) Injecting higher noise magnitudes (σ_{noise} = 399 (0.0, 0.0001, 0.0005, 0.001) leads to a faster decrease in population vector correlation (R_{PV}) from 400 trial 25,000 to 200,000. (C) Agents' reward maximization performance (G) remains fairly stable 401 when the noise magnitude increases. Beyond $\sigma_{noise} = 0.001$, performance becomes highly unsta-402 ble. Black dash indicates the trial at which PV and similarity matrix correlation was measured from. 403 (D) The representation similarity matrix (dot product of population activity from (A)) remains sta-404 ble between trials. (E) With higher noise magnitudes, the similarity matrix correlation (R_{RS}) across 405 trials decreases but at a slower rate than PV correlation. (F) Normalized variance in field parameters 406 $(\theta = \{\alpha, \lambda, \sigma\})$ between trials 25,000 to 200,000 quantifies change in individual place fields spatial selectivity. With no noise (blue) or a larger noise magnitude ($\sigma_{noise} = 0.001$), fields with a larger 407 amplitude experiences a greater change in its parameters. When $\sigma_{noise} \in \{0.0001, 0.00025\}$, we 408 see the opposite trend, where fields with a larger amplitude are more stable than fields with a smaller 409 amplitude, replicating Qin et al. (2023). Shaded area is 95% CI over 10 seeds. 410

(Fig. 3D), even with a higher noise magnitude (Fig. 3E), although the agents are not explicitly optimizing for representational similarity (Qin et al., 2023). Unlike noisy field parameter updates, adding noise to the actor and critic synapses caused the agent's reward maximization behavior, representation similarity correlation and population vector correlation to change at similar rates (Sup. Fig. 7), which is not as consistent with experiment (Sup. Fig. 9 for comparisons to data).

We quantified this drifting behavior at the level of individual neurons by summing the normalized 417 (between [0,1]) variance in each field's parameters $(\sum Var(\tilde{\theta}) = Var(\tilde{\alpha}) + Var(\tilde{\lambda}) + Var(\tilde{\sigma}))$ 418 across learning trials, and comparing this against the mean amplitudes for each field. When no 419 Gaussian noise is added (Fig. 3F), fields with a higher mean amplitude showed a higher variance in 420 its parameters, which is expected since fields with a higher amplitude are more likely to be involved 421 in policy learning. Conversely, with a small Gaussian noise, we see the opposite trend where fields 422 with a smaller mean amplitude showed a higher variance in parameters while fields with a higher 423 mean amplitude were more stable. At smaller noise magnitudes, there is a strong positive correlation 424 between higher amplitude fields and the magnitude of actor and critic weights (Sup. Fig. 8). This 425 suggests that high-amplitude fields are more involved in policy learning and thus need stability, whereas less important fields can alter their spatial selectivity, consistent with Oin et al. (2023). 426

427

378

379

380

382

384

386

387

388

389

390

391

392

393

394

396

428 4.4 PLACE FIELD REORGANIZATION IMPROVES POLICY CONVERGENCE

429
430 As the reward-maximizing model recapitulates experimentally-observed changes in place fields, it
431 is natural to ask what computational advantage these representational changes might offer. To probe the contributions of each field parameter to policy learning, we perform ablation experiments. These



Figure 4: Field reorganization and noisy updates improve target learning. (A) Optimizing all 447 three field parameters, amplitude, width and center of randomly distributed fields allowed agents 448 $(N = 16, \sigma = 0.1)$ to attain the highest cumulative discounted reward (G), while fields with fixed 449 field parameters attained the lowest. (B) Optimizing place field widths (σ), followed by field am-450 plitudes (α) and lastly field centers (λ) caused the biggest decrease in the number of trials needed 451 for policy convergence ($T_{G>45}$, attain a running average of G = 45 over 300 trials). As the number 452 of fields increased, the number of trials needed for policy convergence decreased and the computa-453 tional advantage afforded by field optimization extinguished. (C) Agents need to navigate to a target 454 that changed after 50,000 trials $x_r = \{0.5, 0.0, 0.75, -0.25, 0.5\}$. Without noisy field parameter up-455 dates, agents ($N = 128, \sigma = 0.1$) struggled to learn new targets (blue, $\sigma_{noise} = 0.0$). Field updates with different noise magnitudes influenced the policy convergence speed and maximum cumulative 456 reward for subsequent targets, with $\sigma_{noise} = 0.0005$ (red) demonstrating the highest improvement. 457 Shaded area is 95% CI over 50 seeds. 458

ablations are particularly important due to the parameter degeneracies in the model: one can tradeoff the place field amplitudes and the critic and actor weights.

462 We first considered the task of navigating to a single fixed target. Agents with fixed place fields at-463 tained the lowest navigational performance with cumulative reward G plateauing at G = 33 per trial 464 (Fig. 4A), and showed the slowest policy convergence even as the number of fields increased (Fig. 465 4B). Optimizing place field widths (σ) contributed to the greatest improvement in maximum reward and largest decrease in the number of trials for policy convergence (Fig. 4A-B). Optimizing place 466 field amplitudes (α) contributed to the next most significant improvement (Fig. 4A-B). Interestingly, 467 place field center (λ) optimization did not contribute to a significant improvement in performance, 468 and in fact caused a significant decrease in reward maximization performance and speed of policy 469 convergence when optimized together with the amplitude parameter. Hence, optimizing field widths 470 followed by amplitudes and lastly centers significantly improved agent's reward maximization per-471 formance and increased the speed of policy convergence. However, as the number of place fields 472 increase (Fig. 4B), the computational advantage afforded by place field optimization extinguishes. 473 Nevertheless, optimizing all the parameters in a small number of fields, e.g. 8, leads to a similar rate 474 of policy convergence than with a larger number of randomly initialized fields e.g. 128, which hints 475 that representation flexibility could allow efficient learning in systems with few neurons. 476

We now turn to the influence of noisy fields when learning to navigate to new targets, inspired by 477 Dohare et al. (2024). With the same random field initialization, agents now have to navigate from the 478 same start location to a target that repeatedly changes location. Although all agents learned to navi-479 gate to the first and the second targets equally well, agents without noisy field updates struggled to 480 learn the next three targets, and achieved a lower average cumulative reward (Figure 4C). Increasing 481 the noise magnitude led to a monotonic improvement in new target learning. Some fields coding for 482 the initial reward location shifted to code for the new reward location (Sup. Fig. 3). However, noise magnitudes beyond a threshold ($\sigma_{noise} = 0.001$) caused average cumulative reward to decrease. 483 These results suggests that there is a functional role for noise, especially for new target learning. We 484 see a similar improvement in reward maximization performance with noisy field updates in a 2D 485 arena with an obstacle when we either change the target or the obstacle location (Sup. Fig. 12).

486 5 DISCUSSION

487 488

We present a two-layer navigation model which uses tunable place fields as feature inputs to an actor 489 and a critic for policy learning. The parameters of the place fields and the policy and value function 490 learn to maximize rewards using the temporal difference (TD) error. Our simple reinforcement 491 learning model reproduces three experimentally-observed neural phenomena: (1) the emergence of 492 a high place field density at rewards, (2) enlargement of fields against the trajectory, and (3) drifting fields without influencing task performance. We analyzed the model to understand how the TD 493 error, number of place fields and noise magnitudes influenced place field representations. Lastly, we 494 demonstrate that learning place field representations with noisy field parameters improves reward 495 maximization and the rate of policy convergence when learning single and multiple targets. 496

497 The proposed reinforcement learning model might be a sufficient toy model for theoretical analy-498 sis (Bordelon et al., 2024) while remaining biologically grounded enough to make experimentally testable predictions (Kumar et al., 2024a). For instance, our model gives an alternative normative 499 account for field elongation against the trajectory, which can be contrasted with the successor repre-500 sentation algorithm (Kumar et al., 2024b; Raju et al., 2024). As the dynamics of fields are different in 501 these two models, they could be distinguishable by experiments that track fields over the full course 502 of learning (Fig. 2C-E, Sup. Fig. 6). Furthermore, place field width and amplitude optimization 503 increases maximum cumulative reward and accelerates policy convergence (Fig. 4A-B). 504

Most models that characterized representational drift were not studied under the context of naviga-505 tional policy learning (Pashakhanloo & Koulakov, 2023; Qin et al., 2023; Ratzon et al., 2024). We 506 showed that increasing the noise magnitudes caused different drift regimes (Fig. 3F; Sup. Fig. 9D), 507 and at very high noise levels navigation behavior started to collapse (Fig. 3C, Sup. Fig. 7). Impor-508 tantly, we showed that fields in the noisy regime allowed agents to consistently learn new targets in 509 both 1D (Fig. 4C) and 2D (Sup. Fig. 12A-B) environments, without getting stuck in local minima. 510 The biological origins of adding noise to place field parameters can be attributed to noisy synaptic 511 plasticity mechanisms (Attardo et al., 2015; Kappel et al., 2015; Mongillo et al., 2017). Other mech-512 anisms such as unstable dynamics in downstream networks (Sorscher et al., 2023) and modulatory 513 mechanisms such dopamine fluctuations (Krishnan & Sheffield, 2023) could adaptively control drift 514 rates. A difficult experiment that could directly verify our model is to induce or constrain place 515 field drift rates in animals and determine how this perturbation influences new target learning. How fluctuations in dopamine, stochastic actions and stochastic firing rates within place fields drive drift 516 rates needs to be explored. The current model provides a starting point for this investigation. 517

518 The proposed model is not without limitations. First, we modeled single peaked place fields instead 519 of the complex representations resulting from single "place" cells, which can be multi-field and 520 multi-scale. Nevertheless, the proposed online reinforcement learning framework is general enough 521 to accommodate other models of place cell description (Mainali et al., 2024; Sorscher et al., 2023)) e.g. Sup. Fig. 14, and can be extended to study representation learning in other brain regions e.g. 522 medial entorhinal (Boccara et al., 2019) or posterior parietal (Suhaimi et al., 2022) cortex. Next, 523 place field parameters are optimized by backpropagating the temporal difference error through the 524 actor and critic components (Sup. Fig. 15). Since the motivation was to develop a normative model 525 whose objective was to maximize rewards, this was a reasonable starting point. However, this model 526 must be extended using biologically-plausible learning rules (Lillicrap et al., 2016; Miconi, 2017; 527 Murray, 2019; Nøkland, 2016) before it can in any way be considered mechanistic (Edelmann & 528 Lessmann, 2018; Kempadoo et al., 2016; Krishnan et al., 2022; Lee et al., 2024; Starkweather & 529 Uchida, 2021). Furthermore, place fields reorganize during latent learning in the absence of re-530 wards. While we have only explored reward maximizing objective, extending our model to examine 531 place field reorganization when optimizing for non-reward based objectives (Fang & Stachenfeld, 532 2023; Foster et al., 2000; Low et al., 2018) is crucial. Since our model computes gradients using the objective, this should be feasible. While our computational experiments successfully demonstrated 533 the model's effectiveness in reproducing three disparate phenomena, further work should test its 534 robustness across other reinforcement learning algorithms e.g. policy gradient (Kumar & Pehle-535 van, 2024). Additionally, we need to explore how place field reorganization scales in larger, more 536 complex environments beyond the few 2D environments we considered. 537

538

540	REFERENCES
541	KEI EKENCE

542 543 544	Ryunosuke Amo, Sara Matias, Akihiro Yamanaka, Kenji F Tanaka, Naoshige Uchida, and Mitsuko Watabe-Uchida. A gradual temporal shift of dopamine responses mirrors the progression of temporal difference error in machine learning. <i>Nature neuroscience</i> , 25(8):1082–1092, 2022.
545 546	Angelo Arleo and Wulfram Gerstner. Spatial cognition and neuro-mimetic navigation: a model of hippocampal place cell activity. <i>Biological cybernetics</i> , 83(3):287–299, 2000.
547 548 549	Alessio Attardo, James E Fitzgerald, and Mark J Schnitzer. Impermanence of dendritic spines in live adult ca1 hippocampus. <i>Nature</i> , 523(7562):592–596, 2015.
550 551	Charlotte N Boccara, Michele Nardin, Federico Stella, Joseph O'Neill, and Jozsef Csicsvari. The entorhinal cognitive map is attracted to goals. <i>Science</i> , 363(6434):1443–1447, 2019.
552 553 554	Blake Bordelon, Paul Masset, Henry Kuo, and Cengiz Pehlevan. Loss dynamics of temporal difference reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
555 556 557	Michael A Brown and Patricia E Sharp. Simulation of spatial learning in the morris water maze by a neural network model of the hippocampal formation and nucleus accumbens. <i>Hippocampus</i> , 5 (3):171–188, 1995.
558 559 560	Daniel Bush, Caswell Barry, Daniel Manson, and Neil Burgess. Using grid cells for navigation. <i>Neuron</i> , 87(3):507–520, 2015.
561 562 563	Marco Contreras, Tatiana Pelc, Martin Llofriu, Alfredo Weitzenfeld, and Jean-Marc Fellous. The ventral hippocampus is involved in multi-goal obstacle-rich spatial navigation. <i>Hippocampus</i> , 28 (12):853–866, 2018.
564 565 566	Peter Dayan. Improving generalization for temporal difference learning: The successor representation. <i>Neural computation</i> , 5(4):613–624, 1993.
567 568	Peter Dayan and Laurence F Abbott. <i>Theoretical neuroscience: computational and mathematical modeling of neural systems</i> . MIT press, 2005.
569 570	Mitchell L de Snoo, Adam MP Miller, Adam I Ramsaran, Sheena A Josselyn, and Paul W Frankland. Exercise accelerates place cell representational drift. <i>Current Biology</i> , 33(3):R96–R97, 2023.
572 573 574	Shibhansh Dohare, J Fernando Hernandez-Garcia, Qingfeng Lan, Parash Rahman, A Rupam Mah- mood, and Richard S Sutton. Loss of plasticity in deep continual learning. <i>Nature</i> , 632(8026): 768–774, 2024.
575 576	Can Dong, Antoine D Madar, and Mark EJ Sheffield. Distinct place cell dynamics in ca1 and ca3 encode experience in new environments. <i>Nature communications</i> , 12(1):2977, 2021.
577 578 579 580	David Dupret, Joseph O'neill, Barty Pleydell-Bouverie, and Jozsef Csicsvari. The reorganization and reactivation of hippocampal maps predict spatial memory performance. <i>Nature neuroscience</i> , 13(8):995–1002, 2010.
581 582	Elke Edelmann and Volkmar Lessmann. Dopaminergic innervation and modulation of hippocampal networks. <i>Cell and tissue research</i> , 373:711–727, 2018.
583 584 585 586	Tamir Eliav, Shir R Maimon, Johnatan Aljadeff, Misha Tsodyks, Gily Ginosar, Liora Las, and Nachum Ulanovsky. Multiscale representation of very large environments in the hippocampus of flying bats. <i>Science</i> , 372(6545):eabg4020, 2021.
587 588	Ching Fang and Kimberly L Stachenfeld. Predictive auxiliary objectives in deep rl mimic learning in the brain. <i>arXiv preprint arXiv:2310.06089</i> , 2023.
589 590 591	Ila R Fiete, Yoram Burak, and Ted Brookings. What grid cells convey about rat location. <i>Journal of Neuroscience</i> , 28(27):6858–6871, 2008.
592 593	Stan B Floresco, Christopher L Todd, and Anthony A Grace. Glutamatergic afferents from the hip- pocampus to the nucleus accumbens regulate activity of ventral tegmental area dopamine neurons. <i>Journal of Neuroscience</i> , 21(13):4915–4922, 2001.

594 David J Foster, Richard GM Morris, and Peter Dayan. A model of hippocampally dependent navi-595 gation, using the temporal difference learning rule. *Hippocampus*, 10(1):1–16, 2000. 596 Loren M Frank, Garrett B Stanley, and Emery N Brown. Hippocampal plasticity across multiple 597 days of exposure to novel environments. Journal of Neuroscience, 24(35):7681–7689, 2004. 598 Nicolas Frémaux, Henning Sprekeler, and Wulfram Gerstner. Reinforcement learning using a con-600 tinuous time actor-critic framework with spiking neurons. PLoS computational biology, 9(4): 601 e1003024, 2013. 602 Matthew PH Gardner, Geoffrey Schoenbaum, and Samuel J Gershman. Rethinking dopamine as 603 generalized prediction error. Proceedings of the Royal Society B, 285(1891):20181645, 2018. 604 605 Jeffrey L Gauthier and David W Tank. A dedicated population for reward coding in the hippocam-606 pus. Neuron, 99(1):179–193, 2018. 607 Tom M George, William de Cothi, Kimberly L Stachenfeld, and Caswell Barry. Rapid learning of 608 predictive maps with stdp and theta phase precession. Elife, 12:e80663, 2023. 609 610 Samuel J Gershman. The successor representation: its computational logic and neural substrates. 611 Journal of Neuroscience, 38(33):7193–7200, 2018. 612 Samuel J Gershman and Naoshige Uchida. Believing in dopamine. Nature Reviews Neuroscience, 613 20(11):703–714, 2019. 614 Nitzan Geva, Daniel Deitch, Alon Rubin, and Yaniv Ziv. Time and experience differentially affect 615 distinct aspects of hippocampal representational drift. Neuron, 111(15):2357-2366, 2023. 616 617 Walter G Gonzalez, Hanwen Zhang, Anna Harutyunyan, and Carlos Lois. Persistence of neuronal 618 representations through time and damage in the hippocampus. Science, 365(6455):821-825, 619 2019. 620 Bruce Harland, Marco Contreras, Madeline Souder, and Jean-Marc Fellous. Dorsal ca1 hippocampal 621 place cells form a multi-scale representation of megaspace. Current Biology, 31(10):2178–2190, 622 2021. 623 624 James C. Houk, James L. Adams, and Andrew G. Barto. A Model of How the Basal Ganglia Gen-625 erate and Use Neural Signals That Predict Reinforcement. In Models of Information Processing in the Basal Ganglia. The MIT Press, 11 1994. ISBN 9780262275774. doi: 10.7551/mitpress/ 626 4708.003.0020. URL https://doi.org/10.7551/mitpress/4708.003.0020. 627 628 Daphna Joel, Yael Niv, and Eytan Ruppin. Actor-critic models of the basal ganglia: New anatomical 629 and computational perspectives. Neural networks, 15(4-6):535-547, 2002. 630 Min W Jung, Sidney I Wiener, and Bruce L McNaughton. Comparison of spatial firing character-631 istics of units in dorsal and ventral hippocampus of the rat. Journal of Neuroscience, 14(12): 632 7347-7356, 1994. 633 634 David Kappel, Stefan Habenschuss, Robert Legenstein, and Wolfgang Maass. Network plasticity as 635 bayesian inference. PLoS computational biology, 11(11):e1004485, 2015. 636 Kimberly A Kempadoo, Eugene V Mosharov, Se Joon Choi, David Sulzer, and Eric R Kandel. 637 Dopamine release from the locus coeruleus to the dorsal hippocampus promotes spatial learning 638 and memory. Proceedings of the National Academy of Sciences, 113(51):14835–14840, 2016. 639 640 Clifford G Kentros, Naveen T Agnihotri, Samantha Streater, Robert D Hawkins, and Eric R Kan-641 del. Increased attention to spatial context increases both place field stability and spatial memory. Neuron, 42(2):283-295, 2004. 642 643 Seetha Krishnan and Mark EJ Sheffield. Reward expectation reduces representational drift in the 644 hippocampus. bioRxiv, 2023. 645 Seetha Krishnan, Chad Heer, Chery Cherian, and Mark EJ Sheffield. Reward expectation extinc-646 tion restructures and degrades cal spatial maps through loss of a dopaminergic reward proximity 647

signal. Nature communications, 13(1):6662, 2022.

M Ganesh Kumar and Cengiz Pehlevan. Place fields organize along goal trajectory with reinforce-649 ment learning. Cognitive Computational Neuroscience, 2024. 650 M Ganesh Kumar, Cheston Tan, Camilo Libedinsky, Shih-Cheng Yen, and Andrew YY Tan. A 651 nonlinear hidden layer enables actor-critic agents to learn multiple paired association navigation. 652 Cerebral Cortex, 32(18):3917-3936, 2022. 653 654 M Ganesh Kumar, Shamini Ayyadhury, and Elavazhagan Murugan. Trends innovations challenges 655 in employing interdisciplinary approaches to biomedical sciences. In Translational Research in 656 Biomedical Sciences: Recent Progress and Future Prospects, pp. 287-308. Springer, 2024a. 657 M Ganesh Kumar, Cheston Tan, Camilo Libedinsky, Shih-Cheng Yen, and Andrew Yong-Yi Tan. 658 One-shot learning of paired association navigation with biologically plausible schemas, 2024b. 659 URL https://arxiv.org/abs/2106.03580. 660 661 Jae Sung Lee, John J Briguglio, Jeremy D Cohen, Sandro Romani, and Albert K Lee. The statistical 662 structure of the hippocampal code for space as a function of time, context, and value. Cell, 183 663 (3):620–635, 2020. 664 665 Rachel S Lee, Yotam Sagiv, Ben Engelhard, Ilana B Witten, and Nathaniel D Daw. A featurespecific prediction error model explains dopaminergic heterogeneity. *Nature neuroscience*, 27(8): 666 1574-1586, 2024. 667 668 Timothy P Lillicrap, Daniel Cownden, Douglas B Tweed, and Colin J Akerman. Random synaptic 669 feedback weights support error backpropagation for deep learning. *Nature communications*, 7(1): 670 13276, 2016. 671 John E Lisman and Anthony A Grace. The hippocampal-vta loop: controlling the entry of informa-672 tion into long-term memory. *Neuron*, 46(5):703–713, 2005. 673 674 Ryan J Low, Sam Lewallen, Dmitriy Aronov, Rhino Nevers, and David W Tank. Probing variability 675 in a cognitive map using manifold inference from neural dynamics. *BioRxiv*, pp. 418939, 2018. 676 677 Nischal Mainali, Rava Azeredo da Silveira, and Yoram Burak. Universal statistics of hippocampal 678 place fields across species and dimensionalities. *bioRxiv*, pp. 2024–06, 2024. 679 Emily A Mankin, Fraser T Sparks, Begum Slayyeh, Robert J Sutherland, Stefan Leutgeb, and Jill K 680 Leutgeb. Neuronal code for extended time in the hippocampus. Proceedings of the National 681 Academy of Sciences, 109(47):19462–19467, 2012. 682 683 Paul Masset, Shanshan Qin, and Jacob A Zavatone-Veth. Drifting neuronal representations: Bug or 684 feature? Biological cybernetics, 116(3):253–266, 2022. 685 Mayank R Mehta, Carol A Barnes, and Bruce L McNaughton. Experience-dependent, asymmetric 686 expansion of hippocampal place fields. Proceedings of the National Academy of Sciences, 94(16): 687 8918-8921, 1997. 688 689 Mayank R Mehta, Michael C Quirk, and Matthew A Wilson. Experience-dependent asymmetric 690 shape of hippocampal receptive fields. Neuron, 25(3):707-715, 2000. 691 Ishai Menache, Shie Mannor, and Nahum Shimkin. Basis function adaptation in temporal difference 692 reinforcement learning. Annals of Operations Research, 134(1):215-238, 2005. 693 694 Thomas Miconi. Biologically plausible learning in recurrent neural networks reproduces neural 695 dynamics observed during cognitive tasks. *Elife*, 6:e20899, 2017. 696 697 Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), Proceedings of The 33rd 699 International Conference on Machine Learning, volume 48 of Proceedings of Machine Learning 700 Research, pp. 1928–1937, New York, New York, USA, 20–22 Jun 2016. PMLR. URL https: 701 //proceedings.mlr.press/v48/mniha16.html.

702 703 704 705	Gianluigi Mongillo, Simon Rumpel, and Yonatan Loewenstein. Intrinsic volatility of synaptic con- nections—a challenge to the synaptic trace theory of memory. <i>Current opinion in neurobiology</i> , 46:7–13, 2017.
706 707 708	P Read Montague, Peter Dayan, and Terrence J Sejnowski. A framework for mesencephalic dopamine systems based on predictive hebbian learning. <i>Journal of neuroscience</i> , 16(5):1936–1947, 1996.
709 710 711	Richard GM Morris, Paul Garrud, JNP al Rawlins, and John O'Keefe. Place navigation impaired in rats with hippocampal lesions. <i>Nature</i> , 297(5868):681–683, 1982.
712 713	May-Britt Moser, David C Rowland, and Edvard I Moser. Place cells, grid cells, and memory. <i>Cold Spring Harbor perspectives in biology</i> , 7(2):a021808, 2015.
714 715 716	James M Murray. Local online learning in recurrent networks with random feedback. <i>Elife</i> , 8: e43299, 2019.
717 718 710	Yael Niv. Reinforcement learning in the brain. <i>Journal of Mathematical Psychology</i> , 53(3):139–154, 2009.
719 720 721	Arild Nøkland. Direct feedback alignment provides learning in deep neural networks. Advances in neural information processing systems, 29, 2016.
722 723	J O'Keefe. The hippocampus as a cognitive map, 1978.
724 725	John O'Keefe and Jonathan Dostrovsky. The hippocampus as a spatial map: preliminary evidence from unit activity in the freely-moving rat. <i>Brain research</i> , 1971.
726 727 728 729	Mark G Packard and James L McGaugh. Inactivation of hippocampus or caudate nucleus with lido- caine differentially affects expression of place and response learning. <i>Neurobiology of learning</i> <i>and memory</i> , 65(1):65–72, 1996.
730 731 732	Jon Palacios-Filardo and Jack R Mellor. Neuromodulation of hippocampal long-term synaptic plas- ticity. <i>Current opinion in neurobiology</i> , 54:37–43, 2019.
733 734 735	Farhad Pashakhanloo and Alexei Koulakov. Stochastic gradient descent-induced drift of represen- tation in a two-layer neural network. In <i>International Conference on Machine Learning</i> , pp. 27401–27419. PMLR, 2023.
736 737 738 739	James B Priestley, John C Bowler, Sebi V Rolotti, Stefano Fusi, and Attila Losonczy. Signatures of rapid plasticity in hippocampal cal representations during novel experiences. <i>Neuron</i> , 110(12): 1978–1992, 2022.
740 741 742	Shanshan Qin, Shiva Farashahi, David Lipshutz, Anirvan M Sengupta, Dmitri B Chklovskii, and Cengiz Pehlevan. Coordinated drift of receptive fields in hebbian/anti-hebbian network models during noisy representation learning. <i>Nature Neuroscience</i> , 26(2):339–349, 2023.
743 744 745 746	Rajkumar Vasudeva Raju, J Swaroop Guntupalli, Guangyao Zhou, Carter Wendelken, Miguel Lázaro-Gredilla, and Dileep George. Space is a latent sequence: A theory of the hippocampus. <i>Science Advances</i> , 10(31):eadm8470, 2024.
747 748	Aviv Ratzon, Dori Derdikman, and Omri Barak. Representational drift as a result of implicit regularization. <i>Elife</i> , 12:RP90069, 2024.
749 750 751	John NJ Reynolds, Brian I Hyland, and Jeffery R Wickens. A cellular mechanism of reward-related learning. <i>Nature</i> , 413(6851):67–70, 2001.
752 753	Uri Rokni, Andrew G Richardson, Emilio Bizzi, and H Sebastian Seung. Motor learning with unstable neural representations. <i>Neuron</i> , 54(4):653–666, 2007.
/ 34	Spott I Dusso and Eric I Nostlar. The brain reward aircuitry in mood disorders. Nature reviews

755 Scott J Russo and Eric J Nestler. The brain reward circuitry in mood disorders. *Nature reviews neuroscience*, 14(9):609–625, 2013.

- 756 Fares JP Sayegh, Lionel Mouledous, Catherine Macri, Juliana Pi Macedo, Camille Lejards, Claire Rampon, Laure Verret, and Lionel Dahan. Ventral tegmental area dopamine projections to the 758 hippocampus trigger long-term potentiation and contextual learning. Nature Communications, 15 759 (1):4100, 2024. 760 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-761 dimensional continuous control using generalized advantage estimation. arXiv preprint 762 arXiv:1506.02438, 2015. 763 764 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy 765 optimization algorithms. arXiv preprint arXiv:1707.06347, 2017. 766 Wolfram Schultz, Peter Dayan, and P Read Montague. A neural substrate of prediction and reward. 767 Science, 275(5306):1593-1599, 1997. 768 769 Ben Sorscher, Gabriel C Mel, Samuel A Ocko, Lisa M Giocomo, and Surya Ganguli. A unified 770 theory for the computational and mechanistic origins of grid cells. Neuron, 111(1):121-137, 771 2023. 772 Marielena Sosa, Mark H Plitt, and Lisa M Giocomo. Hippocampal sequences span experience 773 relative to rewards. bioRxiv, 2023. 774 775 Kimberly L Stachenfeld, Matthew M Botvinick, and Samuel J Gershman. The hippocampus as a predictive map. Nature neuroscience, 20(11):1643-1653, 2017. 776 777 Clara Kwon Starkweather and Naoshige Uchida. Dopamine signals as temporal difference errors: 778 recent advances. Current Opinion in Neurobiology, 67:95–105, 2021. 779 RJ Steele and RGM Morris. Delay-dependent impairment of a matching-to-place task with chronic 780 and intrahippocampal infusion of the nmda-antagonist d-ap5. *Hippocampus*, 9(2):118–136, 1999. 781 782 Ahmad Suhaimi, Amos WH Lim, Xin Wei Chia, Chunyue Li, and Hiroshi Makino. Representa-783 tion learning in the artificial and biological neural networks underlying sensorimotor integration. 784 Science Advances, 8(22):eabn0984, 2022. 785 Richard S. Sutton and Andrew G. Barto. Reinforcement learning: An introduction. A Bradford 786 Book, 2018. 787 788 Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient meth-789 ods for reinforcement learning with function approximation. Advances in neural information 790 processing systems, 12, 1999. 791 Jane X Wang, Zeb Kurth-Nelson, Dharshan Kumaran, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, 792 Demis Hassabis, and Matthew Botvinick. Prefrontal cortex as a meta-reinforcement learning 793 system. Nature neuroscience, 21(6):860-868, 2018. 794 Sara Zannone, Zuzanna Brzosko, Ole Paulsen, and Claudia Clopath. Acetylcholine-modulated plas-796 ticity in reward-driven navigation: a computational study. *Scientific reports*, 8(1):9486, 2018. 797 Yaniv Ziv, Laurie D Burns, Eric D Cocker, Elizabeth O Hamel, Kunal K Ghosh, Lacey J Kitch, 798 Abbas El Gamal, and Mark J Schnitzer. Long-term dynamics of ca1 hippocampal place codes. 799 Nature neuroscience, 16(3):264-266, 2013. 800 801 802 803 804 805
- 808
- 809

A DETAILS OF THE PLACE FIELD-BASED NAVIGATION MODEL

A.1 PLACE FIELDS IN 1D AND 2D ENVIRONMENTS

The agent contains N place fields. In a 1D track, each place field is described as

$$\phi_i(x_t) = \alpha_i^2 \exp\left(-\frac{||x_t - \lambda_i||_2^2}{2\sigma_i^2}\right),\tag{8}$$

with α , λ and σ describing the amplitude, center and width, adapted from Foster et al. (2000); Kumar et al. (2022; 2024b). Most of the simulations were initialized with amplitudes $\alpha_i = 0.5$ and widths $\sigma_i = 0.1$, with centers uniformly tiling the environment $\lambda = \{-1, ..., 1\}$. Nevertheless, similar representations emerge for amplitudes drawn from a uniform distribution between [0, 1] and widths uniformly drawn between [0.01, 0.25]. This parameter initialization was used for ablation studies in Fig. 4. In a 2D arena, each place field is described as

$$\phi_i(x_t) = \alpha_i^2 \exp\left[-\frac{1}{2}(x_t - \lambda_i)^\top \Sigma_i^{-1}(x_t - \lambda_i)\right], \qquad (9)$$

where Σ_i is a 2x2 covariance matrix, adapated from Menache et al. (2005). The off-diagonals were initialized as zeros and diagonals initialized to match the variance in the 1D place field description, i.e. $\Sigma_{ii} = 0.1^2$ to ensure field widths are consistent in 1D and 2D.

A.2 REWARD MAXIMIZATION OBJECTIVE (POLICY GRADIENT)

The objective of the model is to learn a policy π parametrized by W^{π} and spatial features ϕ parameterized by θ that maximizes the expected cumulative discounted rewards over trajectories τ in a finite-horizon setting, modeling the trial structure in neuroscience experiments

$$\mathcal{J}^{G} = \mathbb{E}_{\tau \sim \phi_{\theta}, \pi_{W}\pi} \left[\sum_{t=0}^{T} \sum_{k=0}^{T} \gamma^{k} r_{t+1+k} \right] = \mathbb{E} \left[\sum_{t=0}^{T} G_{t} \right], \tag{10}$$

where γ is the discount factor, r_{t+1} is the reward at time step t + 1 after choosing an action g_t at time step t, and the time horizon T is finite with trials ending after a maximum of 100 steps in the 1D track and 300 steps in the 2D arena.

To maximize the cumulative reward objective, we perform gradient ascent on the policy and place field parameters,

$$\theta_{new} = \theta_{old} + \eta_{\theta} \nabla_{\theta} \mathcal{J}^G \quad , \quad W_{new}^{\pi} = W_{old}^{\pi} + \eta \nabla_{W^{\pi}} \mathcal{J}^G \,, \tag{11}$$

where η_{θ} and η are learning rates for θ and W^{π} respectively. The gradients are derived using the log-derivative trick,

$$\nabla_{\theta, W^{\pi}} \mathcal{J}^{G} = \nabla_{\theta, W^{\pi}} \mathbb{E}\left[G(\tau)\right]$$
(12)

$$= \nabla_{\theta, W^{\pi}} \int_{\tau} p(\tau | \theta, W^{\pi}) G(\tau)$$
(13)

$$= \int p(\tau|\theta, W^{\pi}) \nabla_{\theta, W^{\pi}} \log p(\tau|\theta, W^{\pi}) G(\tau)$$
(14)

$$= \mathbb{E}\left[\nabla_{\theta, W^{\pi}} \log p(\tau | \theta, W^{\pi}) G(\tau)\right], \qquad (15)$$

where the trajectory τ describes the state to state transitions. We expand the above using the Markov assumption that the transition to future states depend only on the present state and not on the states preceding it,

$$p(\tau|\theta, W^{\pi}) = p(x_0) \prod_{t=0}^{T} p(x_{t+1}|x_t) \pi(g_t|x_t; \theta, W^{\pi})$$
(16)

$$\log p(\tau|\theta, W^{\pi}) = \log p(x_0) + \sum_{t=0}^{I} \left(\log p(x_{t+1}|x_t) + \log \pi(g_t|x_t; \theta, W^{\pi})\right)$$
(17)

$$\nabla_{\theta, W^{\pi}} \log p(\tau | \theta, W^{\pi}) = \sum_{t=0}^{T} \log \pi(g_t | x_t; \theta, W^{\pi}).$$
(18)

Since the gradients are not dependent on the state transitions, the last line excludes them. Substituting Eq. 18 into Eq. 15 yields

$$\nabla_{\theta, W^{\pi}} \mathcal{J}^{G} = \mathbb{E} \left[\sum_{t=0}^{T} \nabla_{\theta, W^{\pi}} \log \pi(g_{t} | x_{t}; \theta, W^{\pi}) \cdot G_{t} \right],$$
(19)

which completes the full derivation of the policy gradient theorem (Sutton & Barto, 2018; Sutton et al., 1999). The policy gradient objective was used by Kumar & Pehlevan (2024) to optimize the policy and place field parameters. However, this learning signal requires an explicit reward and policy gradient methods are slow to converge as they suffer from high variance due to:

- Monte Carlo sampling: Agents need to sample an entire episode to estimate the expected return $\mathbb{E}_{\tau}[G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + ...]$ before updating the policy. This can introduce significant variance because the estimate is based on a single path through the stochastic environment, which may not be representative of the expected value over many episodes.
 - No Baseline: The basic policy gradient algorithm computes the gradient solely based on the return G from each trajectory. By introducing a baseline (either constant b or dynamically evolving b_t e.g. value function v_t), which estimates the expected return from a given state, the variance of the gradient estimate can be reduced, because now the policy learns which action is better than the previous (concept of using an Advantage A_t instead of rewards).

Value based methods (Sutton & Barto (2018), Chapter 3.5) were introduced to address some of these issues. For instance, instead of sampling returns G_t , value functions V_t learn to estimate the expected returns

$$V_t = \mathbb{E}[G_t], \qquad (20)$$

which can reduce the variance during credit assignment. The combination of policy gradient with value-based methods lead us to the Actor-Critic algorithm.

A.3 ALTERNATIVE REWARD MAXIMIZATION OBJECTIVE (TEMPORAL DIFFERENCE)

The optimal value function V_t reflects the true expected cumulative discounted rewards, hence the policy gradient objective can be rewritten as

$$\mathcal{J}^G = \mathbb{E}\left[\sum_{t=0}^T G_t\right] = \mathbb{E}\left[\sum_{t=0}^T \sum_{k=0}^T \gamma^k r_{t+1+k}\right] = \sum_{t=0}^T V_t, \qquad (21)$$

$$= \mathbb{E}\left[\sum_{t=0}^{T} r_{t+1} + \gamma \sum_{k=0}^{T} \gamma^k r_{t+2+k}\right], \qquad (22)$$

$$\mathcal{J}^G = \mathbb{E}\left[\sum_{t=0}^T r_{t+1} + \gamma G_{t+1}\right] = \mathbb{E}\left[\sum_{t=0}^T r_{t+1} + \gamma V_{t+1}\right].$$
(23)

905 which yields the following self-consistency equation

$$r_{t+1} + \gamma V_{t+1} - V_t \equiv 0, \qquad (24)$$

907 as argued by Frémaux et al. (2013); Sutton & Barto (2018).

Alternatives to policy gradient algorithms propose subtracting a baseline which can be a fixed constant b or a dynamically changing variable b_t . Since we have the value function V_t we can modify the objective to be

$$\mathcal{J}^{A} = \mathbb{E}\left[G_{t} - V_{t}\right] = \mathbb{E}\left[A_{t}\right] = \mathbb{E}\left[\sum_{t=0}^{T} r_{t+1} + \gamma V_{t+1} - V_{t}\right], \qquad (25)$$

which gives us the Advantage function (Mnih et al., 2016; Schulman et al., 2015). This reduces
the variance as the policy has to learn to select actions that gives an advantage over the current value function. We get a learning objective function that is an analogue to maximizing the expected

cumulative discounted returns while subtracting a baseline Eq. 10. 5π

$$\nabla_{\theta, W^{\pi}} \mathcal{J}^{A} = \mathbb{E} \left[\sum_{t=0}^{T} \nabla_{\theta} \log \pi(g_{t} | x_{t}; \theta, W^{\pi}) \cdot A_{t} \right].$$
(26)

However, we have assumed that we are given the optimal value function V_t to critique the actor if it is doing better or worse than before. Instead, we can learn to estimate the value function v_t using a critic by minimizing the Temporal Difference error

$$r_{t+1} + \gamma v_{t+1} - v_t = \delta_t \,. \tag{27}$$

⁹²⁶ The critic can learn to approximate the true value function by minimizing the mean squared error ⁹²⁷ between the true value function V_t and the predicted v_t , or the temporal difference error δ_t ⁹²⁸

$$\mathcal{L}^{v} = \mathbb{E}\left[\sum_{t=0}^{T} \frac{1}{2} \left(V(x_t) - v(x_t; \theta, w^v)\right)^2\right]$$
(28)

$$= \mathbb{E}\left[\sum_{t=0}^{T} \frac{1}{2} \left(r_{t+1} + \gamma V(x_{t+1}) - v(x_t; \theta, w^v)\right)^2\right].$$
 (29)

Since we do not have the optimal value function V_t , we can approximate it by bootstrapping the estimated value function v_t and ensuring that we do not take gradients with respect to the time shifted value estimate $v(x_{t+1})$

$$\mathcal{L}^{TD} = \mathbb{E}\left[\sum_{t=0}^{T} \frac{1}{2} \left(r_{t+1} + \gamma v(x_{t+1}) - v(x_t; \theta, w^v)\right)^2\right]$$
(30)

$$= \mathbb{E}\left[\sum_{t=0}^{T} \frac{1}{2} \delta_t^2(\theta, w^v)\right] \,. \tag{31}$$

We minimize the temporal difference error using gradient descent for the critic to estimate the value function

$$\nabla_{\theta, w^{v}} \mathcal{L}^{TD} = \frac{\partial \mathcal{L}^{TD}}{\partial \delta} \cdot \frac{\partial \delta}{\partial v} \cdot \nabla_{\theta, w^{v}} v(\theta, w^{v}), \qquad (32)$$

$$= \mathbb{E}\left[\sum_{t=0}^{T} \delta_t \cdot (-1) \cdot \nabla_{\theta, w^v} v(x_t; \theta, w^v)\right], \qquad (33)$$

$$= \mathbb{E}\left[\sum_{t=0}^{T} -\nabla_{\theta^{v}} v(x_{t}; \theta, w^{v}) \cdot \delta_{t}\right].$$
(34)

Notice the additional negative sign that pops out when you take the derivative of δ only with respect to v_t

=

$$\frac{\partial \delta}{\partial v} = \frac{\partial (r_{t+1} + \gamma v_{t+1} - v_t)}{\partial v_t} = -1, \qquad (35)$$

since r_{t+1} and v_{t+1} are treated as constants, we do not take their derivatives. Since we do not have the optimal value function V_t but a biased estimate v_t , we can use the temporal difference error as our reward maximization objective

$$\mathcal{J}^{TD} = \mathbb{E}\left[\sum_{t=0}^{T} r_{t+1} + \gamma v_{t+1} - v_t\right] = \mathbb{E}\left[\sum_{t=0}^{T} \delta_t\right].$$
(36)

As the value estimation becomes closer to the optimal value $v_t \to V_t$, this objective becomes similar to the advantage objective $\mathcal{J}^{TD} \to \mathcal{J}^A$. Note that we are not directly maximizing the TD error during policy learning. Rather, we want to optimize the policy π and place field parameters θ by gradient ascent, using the biased estimate of the advantage function

969
970
971
$$\nabla_{\theta,W^{\pi}} \mathcal{J}^{TD} = \mathbb{E} \left[\sum_{t=0}^{T} \nabla_{\theta,W^{\pi}} \log \pi(g_t | x_t; \theta, W^{\pi}) \cdot \delta_t \right].$$
(37)

An alternative interpretation is that during policy learning, the agent learns a policy to maximize the difference between the actual reward and the estimated value

COMBINED REWARD MAXIMIZATION OBJECTIVE FOR PLACE FIELD PARAMETERS A.4

In our model (Fig. 1A), actor W^{π} and critic w^{v} weights are optimized separately, while the place field parameters θ overlap. The actor uses gradient ascent for Eq. 26, and the critic employs gradient descent for Eq. 34. Since we have a single population of place fields, we optimize these parameters to support both objectives. Thus, we derive a combined objective function to update W^{π} , w^{v} , and θ in a single gradient pass

$$\nabla_{W^{\pi},w^{\nu},\theta}\mathcal{J} = \nabla_{W^{\pi},w^{\nu},\theta}\mathcal{J}^{TD} - \nabla_{W^{\pi},w^{\nu},\theta}\mathcal{L}^{TD}$$
(38)

$$= \mathbb{E}\left[\sum_{t=0}^{T} \nabla_{W^{\pi}, w^{\upsilon}, \theta} \log \pi(g_t | x_t; W^{\pi}, \theta) \delta_t\right] - \mathbb{E}\left[\sum_{t=0}^{T} -\nabla_{W^{\pi}, w^{\upsilon}, \theta} v(x_t; w^{\upsilon}, \theta) \delta_t\right],$$
(39)

$$= \mathbb{E}\left[\sum_{t=0}^{T} \nabla_{W^{\pi}, w^{v}, \theta} \log \pi(g_{t}|x_{t}; W^{\pi}, \theta) \delta_{t} + \nabla_{W^{\pi}, w^{v}, \theta} v(x_{t}; w^{v}, \theta) \delta_{t}\right],$$
(40)

$$= \mathbb{E}\left[\sum_{t=0}^{T} \left(\nabla_{W^{\pi}, w^{v}, \theta} \log \pi(g_{t}|x_{t}; W^{\pi}, \theta) + \nabla_{W^{\pi}, w^{v}, \theta} v(x_{t}; w^{v}, \theta)\right) \delta_{t}\right].$$
(41)

where $\nabla_{w^v} \mathcal{J}^{TD} = 0$ and $\nabla_{W^{\pi}} \mathcal{L}^{TD} = 0$ since the respective objectives are not parameterized by w^{v} and W^{π} respectively. This means that W^{π} is tuned to maximize \mathcal{J}^{TD} , w^{v} is tuned to minimize \mathcal{L}^{TD} and θ is tuned to balance both the objectives.

Since most optimizers e.g. in Tensorflow, PyTorch perform gradient descent, not ascent, we can minimize the negative policy gradient Eq. 26, which is equivalent to the negative log likelihood

$$\begin{aligned}
& 1001 \\
& 1002 \\
& \nabla_{W^{\pi},w^{v},\theta}\mathcal{L} = -\nabla_{W^{\pi},w^{v},\theta}\mathcal{J}^{TD} + \nabla_{W^{\pi},w^{v},\theta}\mathcal{L}^{TD} \\
& (42)
\end{aligned}$$

$$= -\mathbb{E}\left[\sum_{t=0}^{T} \nabla_{W^{\pi}, w^{v}, \theta} \log \pi(a_{t}|x_{t}; W^{\pi}, \theta) \cdot \delta_{t}\right] + \mathbb{E}\left[\sum_{t=0}^{T} -\nabla_{W^{\pi}, w^{v}, \theta} \tilde{v}(x_{t}; w^{v}, \theta) \cdot \delta_{t}\right]$$

$$(43)$$

$$= \mathbb{E}\left[\sum_{t=0}^{T} \nabla_{W^{\pi}, w^{\upsilon}, \theta} - \log \pi(a_t | x_t; W^{\pi}, \theta) \cdot \delta_t\right] + \mathbb{E}\left[\sum_{t=0}^{T} - \nabla_{W^{\pi}, w^{\upsilon}, \theta} \tilde{v}(x_t; w^{\upsilon}, \theta) \cdot \delta_t\right]$$

$$(44)$$

$$= \nabla_{W^{\pi}, w^{v}, \theta} \mathcal{L}_{\pi}^{TD} + \nabla_{W^{\pi}, w^{v}, \theta} \mathcal{L}_{v}^{TD} \,. \tag{45}$$

which is the same update rule used in Mnih et al. (2016); Wang et al. (2018) to train the actor and critic separately while the feature parameters are trained jointly.

It is also possible to initialize two separate populations of place fields, each for the actor and critic. Alternatively, we only optimize place field parameters using the actor's objective while the critic uses the spatial features to learn the value function. The converse is also possible where the place field parameters and critic weights are optimized to minimize the TD error while the actor learns a policy without optimizing the spatial representations, as we did in the perturbative approximation (App. B). From numerical experiments, optimizing place field parameters using both the actor and critic objectives allowed the agent to achieve the fastest policy convergence and highest cumulative reward performance (Sup. Fig. 15).

A.5 ONLINE UPDATE OF PLACE FIELD AND ACTOR-CRITIC PARAMETERS

Now, we derive an online implementation of Eq. 6 which is the same as Eq. 41, so that all parameters are updated at every time step. Extending from Foster et al. (2000); Kumar et al. (2022), the actor and critic weights are updated according to the gradients

$$\Delta w_i^v(t+1) = \eta \delta_t \phi_i(x_t) \quad , \quad \Delta W_{ji}^\pi(t+1) = \eta \delta_t \phi_i(x_t) \tilde{g}_{t,j}^\top , \tag{46}$$

where $\tilde{g}_{t,j} = g_t - P$ and $\eta = 0.01$. The gradient updates for place field parameters follow

$$\Delta \theta_i(t+1) = \eta_\theta \delta_t \left(w_i^v(t) + W_{ji}^\pi(t) \cdot \tilde{g}_{t,j} \right) \nabla_\theta \phi_i(x_t; \theta_i) , \qquad (47)$$

where we use a significantly smaller learning rate $\eta_{\theta} = 0.0001$ so that the spatial representation evolves in a stable manner. Specifically, each field parameter is updated according to

$$\delta_{i,t}^{bp} = \delta_t \left(w_i^v(t) + W_{ji}^\pi(t) \cdot \tilde{g}_{t,j} \right) , \qquad (48)$$

$$\Delta \alpha_{i,t} = \eta_{\alpha} \cdot \delta_{i,t}^{bp} \cdot \phi_i(x_t) \cdot \left(\frac{2}{\alpha_i}\right) , \qquad (49)$$

$$\Delta \lambda_{i,t} = \eta_{\lambda} \cdot \delta_{i,t}^{bp} \cdot \phi_i(x_t) \cdot \left(\frac{x_t - \lambda_i}{\sigma_i^2}\right) \,, \tag{50}$$

1039 where $\delta_{i,t}^{bp}$ is the TD error gradient that has been backpropagated through the actor and critic weights. 1040 Using just the $w_i^v(t)$ or W_{ji}^{π} weights alone to backpropagate the TD error influences the represen-1041 tation learned by the place field population and ultimately the navigation performance (Sup. Fig. 1042 15).

There are two ways to optimize the place field width parameter. The first and straightforward method is to update the width parameter according to

$$\Delta \sigma_{i,k,t} = \eta_{\sigma} \cdot \delta_{i,t}^{bp} \cdot \phi_{i,k}(x_t) \cdot \left(\frac{(x_t - \lambda_i)^2}{\sigma_{i,k}^3}\right),$$
(51)

where k = 1 in a 1D place field. In a 2D place field with k = 2, we can update the diagonal elements in the 2D matrix while keeping the off-diagonals to zeros as in Menache et al. (2005). However, fields will only elongate along each axis. Instead, in our simulations, we optimized the off-diagonals using the same gradient flow equations. However, we needed to include additional constraints so that each place field's covariance matrix remains 1) symmetric, 2) bounded, and 3)positive semi-definite to perform matrix inversion. Specifically, the covariance matrix was bounded between $[10^{-5}, 0.5]$ to prevent exploding widths and gradients.

¹⁰⁸⁰ B DERIVATION FOR PERTURBATIVE EXPANSION

The dynamics of place field parameters are nonlinear and difficult to characterize analytically. To gain some analytical tractability, we impose a strong separation of timescales between policy learning updates and place field parameter updates. To do so, we set the learning rates for the actor-critic η to be much larger than the learning rates for the place field parameters $\eta_{\alpha}, \eta_{\lambda}, \eta_{\sigma} \ll \eta$. In simulations, we use $\eta = 0.01$ and $\eta_{\theta} = 0.0001$.

¹⁰⁸⁷ The critic estimates the value as

1089 1090

1093 1094

1099

1100

1104

1115

1116

1118

1120

1123 1124

$$v(x_t) = \sum_{i=1}^{N} w_i \phi_i(x_t, \boldsymbol{\theta}_i), \qquad (52)$$

where $\theta_i = (\alpha_i, \lambda_i, \sigma_i)$ are neuron specific parameters (amplitude, mean, and bandwidth respectively). We write w^v as w for clarity. To start with let's just consider

$$\phi_i(x_t, \boldsymbol{\theta}_i) = \alpha_i^2 \exp\left(-\frac{1}{2\sigma_i^2}(x_t - \lambda_i)^2\right).$$
(53)

We consider a TD based update, which in the gradient flow (infinitesimal learning rate) limit can be approximated as

$$\frac{d}{dt}\boldsymbol{w}(t) = \boldsymbol{M}(t)(\boldsymbol{w}^{V} - \boldsymbol{w}(t)), \qquad (54)$$

$$\frac{d}{dt}\boldsymbol{\theta}_i(t) = \epsilon \, w_i(t) \mathbb{E}_{x_t} \nabla_{\boldsymbol{\theta}_i} \phi_i(x_t, \boldsymbol{\theta}_i) \delta_t \,, \tag{55}$$

The key assumption we make is that the dimensionless ratio of learning rates, ϵ is perturbatively small $\eta_{\theta} \ll 1$ (50)

$$\epsilon = \frac{\eta_{\theta}}{\eta} \ll 1,\tag{56}$$

1105 where η_{θ} is the learning rate for the place field parameters θ_i and η is the learning rate for the 1106 actor-critic. The matrix $M(t) = \Sigma(t) - \gamma \Sigma_+(t)$ where $\Sigma = \langle \psi(x_t)\psi(x_t)\rangle$ and $\Sigma_+(t) =$ 1107 $\langle \psi(x_t)\psi(x_{t+1})^{\top}\rangle$ depends on the equal time and time-step shifted correlations of features. The 1108 vector $\boldsymbol{w}^V = \boldsymbol{M}^{-1}\boldsymbol{\Sigma}\boldsymbol{w}_R$ where $\boldsymbol{w}_R \cdot \boldsymbol{\psi}(x) = R(x)$. We investigate a simple perturbation series.

1109 $w(t) = w_0(t) + \epsilon w_1(t) + \epsilon^2 w_2(t) + \dots$ 1110 $2(t) = 2(t) + \epsilon w_1(t) + \epsilon^2 w_2(t) + \dots$

$$\boldsymbol{\theta}(t) = \boldsymbol{\theta}_0(t) + \epsilon \boldsymbol{\theta}_1(t) + \epsilon^2 \boldsymbol{\theta}_2(t) + \dots$$
(57)

and examine the dynamics up to first order in ϵ . We will show that this recovers many qualitative features of the observed representational updates.

1114 The leading zeroth order dynamics are

$$\frac{d}{dt}\boldsymbol{\theta}_0(t) = 0, \ \frac{d}{dt}\boldsymbol{w}_0(t) = \boldsymbol{M}_0(\boldsymbol{w}_V - \boldsymbol{w}_0(t)),$$
(58)

1117 where $M_0 = \Sigma(0) - \gamma \Sigma_+(0)$ is the initial feature covariance under the initial policy.

1119 B.1 PLACE FIELD AMPLITUDE

We start by asserting a separation of timescales between training readout weights and feature parameters during a simple TD learning setup

$$\frac{d}{dt}w_i(t) = \sum_j M_{ij}(w_j^V - w_j), \qquad (59)$$

$$\frac{d}{dt}\alpha_i(t) = \epsilon \, \frac{2}{\alpha_i(t)} w_i \sum_j M_{ij}(w_j^V - w_j) \,, \tag{60}$$

(62)

The zero-th order solution to Eq. 54 is

1131
$$\Delta \boldsymbol{w}_0(t) \equiv \boldsymbol{w}_V - \boldsymbol{w}_0(t) = \exp\left(-\boldsymbol{M}t\right) \boldsymbol{w}_V, \qquad (61)$$

1132
$$w_0(t) = [I - \exp(-Mt)]w_V,$$
1133

which can be substituted in to get the first order correction to the dynamics for θ

1125 1126 1127

1136

1140 1141 1142

1145

1148

1150

$$\frac{d}{dt}\boldsymbol{\alpha}_{1}(t) = 2\boldsymbol{\alpha}_{0}^{-1} \odot [\boldsymbol{I} - \exp\left(-\boldsymbol{M}t\right)] \boldsymbol{w}_{V} \odot \boldsymbol{M} \exp\left(-\boldsymbol{M}t\right) \boldsymbol{w}_{V}.$$
(63)

¹¹³⁷ Under the condition that $\alpha_0 = 1$ and $M = M^{\top}$ we can work out an exact expression in terms of the eigendecomposition $M = \sum_k \lambda_k u_k u_k^{\top}$

$$\boldsymbol{\alpha}_{1}(t) = 2\sum_{k\ell} (\boldsymbol{w}_{V} \cdot \boldsymbol{u}_{k}) (\boldsymbol{u}_{\ell} \cdot \boldsymbol{w}_{V}) (\boldsymbol{u}_{k} \odot \boldsymbol{u}_{\ell}) \left[(1 - e^{-\lambda_{k}t}) - \frac{\lambda_{k}}{\lambda_{k} + \lambda_{\ell}} (1 - e^{-(\lambda_{k} + \lambda_{\ell})t}) \right], \quad (64)$$

we can approximate this at late times as

$$\lim_{t \to \infty} \boldsymbol{\alpha}_1(t) \approx 2\boldsymbol{w}_V \odot \boldsymbol{w}_V. \tag{65}$$

1146 As $t \to \infty$ we can approximate this as $\lim_{t\to\infty} \theta(t) \approx 2(w_V)^2$. This indicates that neurons which 1147 are heavily involved in the reproduction of the value function are upweighted in their amplitude.

1149 B.2 FIELD CENTER

1151 Based on the place field center update equation and rewriting the terms as above,

1158 1159

1169 1170 1171

1174 1175 $\frac{d}{dt}\lambda_i(t) \approx \epsilon \, \frac{x_t - \lambda_i}{\sigma_i^2} \, w_i \phi_i(x) \sum_j \phi_j(x) (w_j^V - w_j) \,. \tag{66}$

We need to compute an average over spatial positions. We approximate the space position early in training as a Gaussian with mean s_0 and variance σ_x^2

$$\left\langle \frac{(x_t - \lambda_i)}{\sigma^2} \phi_i(x) \phi_j(x) \right\rangle \approx \frac{\mu_{ij} - \lambda_i}{\sigma^2} M_{ij} ,$$
 (67)

where $\mu_{ij} = \left(\frac{2}{\sigma^2} + \frac{1}{\sigma_x^2}\right)^{-1} \left(\frac{1}{\sigma^2}(\lambda_i + \lambda_j) + \frac{1}{\sigma_x^2}\bar{\mu}_x\right)$ is the mean value of x obtained by the above Gaussian integral under the approximation that $p(x) \sim \mathcal{N}(\bar{\mu}_x, \sigma_x^2)$. Approximating λ_j as the mean position of the tuning curves $\bar{\lambda}$ we obtain the following prediction

$$\begin{aligned} \mathbf{\lambda}(t) - \mathbf{\lambda}(0) &\approx \epsilon \mathbf{w}^{V} \odot \left[\left(\frac{2}{\sigma^{2}} + \frac{1}{\sigma_{x}^{2}} \right)^{-1} \left(\frac{1}{\sigma^{2}} (\mathbf{\lambda}(0) + \bar{\lambda}) + \frac{1}{\sigma_{x}^{2}} \bar{\mu}_{x} \right) - \mathbf{\lambda}(0) \right] \odot \left[\mathbf{I} - \exp\left(-\mathbf{M}t \right) \right] \mathbf{w}^{V} \\ \mathbf{M}^{V} \\ \end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\begin{aligned} \mathbf{\lambda}(t) - \mathbf{\lambda}(0) &\approx \epsilon \mathbf{w}^{V} \odot \left[\left(\frac{2}{\sigma^{2}} + \frac{1}{\sigma_{x}^{2}} \right)^{-1} \left(\frac{1}{\sigma^{2}} (\mathbf{\lambda}(0) + \bar{\lambda}) + \frac{1}{\sigma_{x}^{2}} \bar{\mu}_{x} \right) - \mathbf{\lambda}(0) \right] \odot \left[\mathbf{I} - \exp\left(-\mathbf{M}t \right) \right] \mathbf{w}^{V} \\ \end{aligned}$$

$$\end{aligned}$$

1168 Following the solution in Eq. 62, we can approximate this at late times as

$$\lim_{t \to \infty} \boldsymbol{\lambda}(t) - \boldsymbol{\lambda}(0) \approx \epsilon \boldsymbol{w}^{V} \odot \left[\left(\frac{2}{\boldsymbol{\sigma}^{2}} + \frac{1}{\sigma_{x}^{2}} \right)^{-1} \left(\frac{1}{\boldsymbol{\sigma}^{2}} (\boldsymbol{\lambda}(0) + \bar{\boldsymbol{\lambda}}) + \frac{1}{\sigma_{x}^{2}} \bar{\mu}_{x} \right) - \boldsymbol{\lambda}(0) \right] \odot \boldsymbol{w}^{V}.$$
(69)

Hence, in addition to the value of a location, three additional factors influence each field's displacement.

$$\lambda_i(t) - \lambda_i(0) \approx \frac{\eta_\lambda}{\eta} \left(\frac{2}{\sigma_i^2} + \frac{1}{\sigma_x^2}\right)^{-1} \left[\frac{\bar{\lambda} - \lambda_i(0)}{\sigma_i^2} + \frac{\bar{\mu}_x - \lambda_i(0)}{\sigma_x^2}\right] w_{v,i}^2(t) , \ \eta_\lambda \ll \eta , \tag{70}$$

where $\bar{\lambda}$ is the agent's expected location sampled from its policy, $\bar{\mu}_x = -0.75$ is the starting location and σ_x is the estimated spread of the trajectory. This analysis suggests that fields will be influenced by both the start location and the location where the agent spends a higher proportion of time at. In later learning phases, this will be the reward location $\bar{\lambda} = 0.5$. Consequently, only the fields near the reward location will shift towards the reward, while the rest of the fields will move towards the start location. We illustrate this perturbative approximation at early and late times of training in Figure 5. The theory is quite accurate early in training, but fails at sufficiently long training time.

- 1183
- 1184
- 1185
- 1186



Figure 5: Difference in early versus late time perturbative approximation. Blue scatter points shows the magnitude and direction of change in (N = 256) field center position compared to the position at which the fields were initialized $(\lambda_i(T) - \lambda_i(0))$. (A) In early time, the perturbative expansion is a good fit to the field center displacement, and captures the shift in fields towards the reward location $x_r = 0.5$ (red) (B) As learning proceeds, the approximation begins to break down for fields further from the reward location. Free parameters were fit with $\overline{\lambda} = 0.535$ and $\sigma_x = 0.45$.

C DETAILS FOR THE SUCCESSOR REPRESENTATION AGENT

The generalized temporal difference error is given by S^{SP}_{R}

$$\delta_{t,i}^{SR} = \phi_j(x_t) + \gamma \psi_j^{\pi}(x_{t+1}) - \psi_j^{\pi}(x_t), \qquad (71)$$

with M_i representing the predicted successor representation and $\phi(x)$ representing the initialized place field representation that is not optimized.

$$\psi_i^{\pi}(x_t) = \sum_{i}^{N} [U_{ji}]_+ \phi_i(x_t) , \qquad (72)$$

The successor representation is computed using a summation of the place fields with a learned matrix U that is positively rectified. The rectification is necessary to have a non-negative representation.

$$\Delta U_t = \phi_i(x_t) \cdot \delta_{t,i}^{\top}, \qquad (73)$$

The matrix U is initialized as an identity matrix and is updated using a two-factor rule using the TD error as in Gardner et al. (2018).

1220 1221

1222

1225

1217

1204 1205

1206

1209

1212

1213

D DETAILS FOR NOISY FIELD UPDATES

To induce drift, we independently introduced noise to field amplitudes, centers and width, as well as the synapses to the actor and critic ($\theta \in \{\alpha, \lambda, \sigma, w^v, W^\pi\}$).

$$_{t+1} = \theta_t + \xi_t \,, \tag{74}$$

1226 where the noise term ξ_t are independent Gaussian noises with zero mean and magnitude $\sigma_{noise} \in \{10^{-6}, 10^{-1}\}$. We performed a noise sweep to determine how increasing the noise magnitude affected the agent's reward maximization behavior, population vector correlation and representation similarity. Refer to Sup. Fig. 7.

1230 1231

- 1232
- 1233
- 1234
- 123
- 1236 1237
- 1029
- 1239
- 1240



Supplementary Figure 1: Influence of place field parameter optimization for a single seed. Ex-1273 ample change in individual field's spatial selectivity ($\phi(x)$, colored), mean firing activity at a location 1274 $(\sum_{i=1}^{N} \phi_i(x))$, field density which is the number of Center of Mass (COM) in a location after smooth-1275 ing with a Gaussian kernel density estimate (gKDE) (gKDE(COM)), blue) and, the frequency of 1276 being in a location $(p_{BM}(x))$, when optimizing different combinations of field parameters $(\alpha, \lambda, \sigma)$ 1277 during reward maximization (RM). The location in which the highest value for mean firing activity, 1278 field density and frequency is attained is indicated by a red, blue and black vertical dash line re-1279 spectively. Optimizing a (A) small number (N = 16) and (B) large number of place fields yields a 1280 similar high mean firing rate at the reward location followed by the start location. However, the field 1281 density evolves differently when in the low field regime, (A) a high density emerges at the reward 1282 location in the early stages of learning, but it shifts to the start location at later stages of learning. (B) In the high field regime, a high field density at the reward location remains stable throughout 1283 learning. Note that COM changes only when the place field centers are optimized ($\Delta\lambda$). Distribu-1284 tion is shown for a single seed run for a homogeneous place field population that has been initialized 1285 by with equal spacing between field centers ($\lambda \in [-1, 1]$), equal amplitude ($\alpha = 0.5$) and width 1286 $(\sigma = 0.01).$ 1287

1289

1290

1291

1292

1293

1294

1295



Supplementary Figure 2: Average change in field density and mean firing rate for different number of place fields. Vertical blue and red dash lines indicate the location with the highest density and mean firing rate, with the legend indicating the location (x). (A) Homogeneous place field distribution was initialized with field parameters similar to Sup. Fig. 1, equal spacing between field centers ($\lambda \in [-1,1]$), equal amplitude ($\alpha = 0.5$) and equal width ($\sigma = 0.01$). (B) All place field parameters center (λ), amplitude (α), and width (σ) were initialized by sampling from a uniform distribution between [-1, 1], [0, 1], $[10^{-5}, 0.1]$ respectively to model heterogeneous place field population. Learning rates for the place field parameters and actor-critic were $n_{\theta} = 0.0001$ and n = 0.01 respectively. Shaded area is 95% CI over 50 different seeds.



Supplementary Figure 3: A small proportion of reward-encoding place fields shift to the new 1381 **reward location.** Agents with N = 256 place fields and Gaussian noise injected to field parameters 1382 $(\sigma_{noise} = 0.0001)$ were trained to navigate to a reward location at $x_r = 0.75$ for 50,000 trials, 1383 thereafter the reward location was shifted to $x_r = -0.2$ for the next 50,000 trials. (A) Place field 1384 density at the start of learning was uniformly distributed (left) and increased near the first reward 1385 location at the end of the first 50,000 trials (center). After the shift in reward location, a high density 1386 of fields emerged at the new reward location (right). The black line shows the learned policy, where 1387 a velocity of 0.1(-0.1) indicates moving right (left). Agents learn to navigate to the reward location, 1388 both before and after the shift. (B) Example distribution of individual place fields before learning 1389 (left), before the shift (center) and after the shift (right). All place field parameters λ , α , and σ were initialized by sampling from a uniform distribution between [-1, 1], [0, 1], $[10^{-5}, 0.1]$ respectively 1390 to model heterogeneous place field population. Notice the backward shift of some place fields 1391 that were at the initial reward location to the new reward location. (C) About 2.6% of the place 1392 fields coding for the initial reward at $x_r = 0.75$ (green dots) shifted to the new reward location 1393 at $x_r = -0.2$ (about 19 of the 734 green dots are within the blue circle). Other place fields at 1394 $x_r = -0.2$ increased their firing rate to encode the new reward location. We see a large number of 1395 fields shifting backward, though not entirely to the new reward location. Shaded area shows 95% 1396 CI for 10 seeds of agents with 256 place fields each. Black and green dots show a total 2560 place fields for all 10 agents. 1398

- 1399
- 1400
- 1401
- 1402
- 1403





Supplementary Figure 4: Weak feature learning with large number of place fields. Critic w_i^v and actor W_{π}^{π} weights were initialized by sampling from a random normal distribution $\mathcal{N}(0, 10^{-5})$, despite the number of place fields N, similar to Foster et al. (2000); Frémaux et al. (2013); Kumar et al. (2022); Zannone et al. (2018). (A) Homogeneous place field population: Place field parameters were initialized with equal spacing between field centers ($\lambda \in [-1, 1]$), equal amplitude ($\alpha = 0.5$) and equal width ($\sigma = 0.01$). (B) Heterogeneous place field population: All place field parameters center (λ), amplitude (α), and width (σ) were initialized by sampling from a uniform distribution between [-1, 1], [0, 1], $[10^{-5}, 0.1]$ respectively. (A-B) The sum of the L2 norm for each place field's center λ , amplitude α and width σ between its initialized and final value decreases as the number of fields available increases. Hence, as the number of fields increases, the change in each place field's parameter becomes smaller. This suggests a weak feature learning regime with large N. (C) Similar to Fig. 1D. Density at the reward location $d(x_r)$ compared to non-reward location d(x') decreases with a higher number of fields. (D) The mean firing rate at the reward location $\sum \phi(x_r)$ compared to non-reward location $\sum \phi(x')$ decreases with a higher number of fields. (C-D) Density and mean firing rate at the reward location are proportional to the reward magnitude (R_{max}) , and inversely proportional to the size of the reward location (R_{size}) . Error bars show 95% CI over 50 different seeds.



Supplementary Figure 5: SR agent architecture and field dynamics. (A) Successor Representation 1490 (SR) agent architecture to learn a navigational policy and the SR place fields. Only the synapses from 1491 the initialized place field (ϕ_{fixed}) to the actor (red) and critic (green), and the synapses (U) to the 1492 SR fields (ψ) were plastic. Refer to App. C for implementation details. (B) Difference in reward 1493 maximization behavior between SR and RM agent, contributing to the dip in correlation between 1494 the proportion of time spent in a location by both agents in Fig. 2D black line. (C) Average change 1495 for 16 place fields' size (firing rate greater than 10^{-3} in the track) (left) and center of mass (right) 1496 when SR and RM agents navigate in a 1D track with the absolute change reflected in the left and 1497 right y axis. Shaded area shows 95% CI over 10 seeds. (D) Spatial representation similarity matrix 1498 for SR (top row) and RM (bottom row) agents in a 1D track is visualized by taking the dot product of the place field activity at each location. (E) Change in individual place field's spatial selectivity 1499 (colored), mean firing rate (red) and frequency of being in a location (black) when fields are learned 1500 using the Successor Representation (Top row) and Reward Maximizing objective (Bottom row). 1501 Top panels T=1000, 10,000 and 50,000 were selected for SR and bottom panels T=1000, 3000 and 1502 50,000 were selected for RM in the paper due to space constraints. 1503

- 1505
- 1506
- 1507
- 1508
- 1509
- 1510
- 1511



Supplementary Figure 6: Field elongation in 2D arena. (A-B) 2D Place field distortion dynam-1556 ics by SR (A) and RM (B) agents as learning proceeds. Numbers in yellow on the obstacle in-1557 dicates (Field ID)-(Maximum firing rate). (C) Average change in 256 field sizes (left) and cen-1558 ter of masses (right) for SR and RM agents navigating in a 2D arena. Shaded area shows 95% 1559 CI over the 256 fields. Note that agent start randomly from three different locations $x_{start} \in$ 1560 $\{(-0.75, -0.75), (-0.75, 0.75), (0.75, 0.75)\}$ to navigate to the target at $x_r = (0.75, -0.75)$. The 1561 change in field COM shows the average change in center of mass with respect to each starting location. Hence, the averaged backward shift in center of mass might not be extensive. Refer to Fig. 1C 1562 for change in goal representation. 1563



Supplementary Figure 7: Noise amplitude monotonically influences population vector correla-tion and agent performance. Adding Gaussian noise with increasing magnitude $[5x10^{-7}, 10^1]$ either in field parameters $(\alpha, \lambda, \sigma)$ or Actor-Critic (W_{π}, w_{ν}) influences the variance in Population Vector Correlation (R_{PV} , blue), Spatial Representation Similarity which is the dot product of field activity (R_{RS} , orange) and cumulative discounted reward (G, green). Low variance of R_{PV} and R_{RS} indicates high correlation as learning progresses. Low variance in G indicates stable perfor-mance. When G increases before decreasing as the noise amplitude increases, agent's navigation performance collapsed and the agent achieves 0 reward with low variance. A high ratio of variance in population vector correlation and reward maximization behavior $(R_{PV}/G, \text{ red})$ indicates that there is an optimal noise amplitude which causes high variance in population vector correlation (low PV correlation) while demonstrating stable performance. A similar analysis can be performed using representational similarity (R_{PV}/R_{RS}) , purple) to determine the optimal noise amplitude for high variance in population vector correlation but stable representation similarity as seen in Qin et al. (2023). Note that our agents are only optimizing for navigation behavior instead of representation similarity.



Supplementary Figure 8: Influence of noisy fields on agent performance and field representa-1652 tion. (A) Reward maximization performance variability increases when noise magnitude increases. 1653 (B) With no noise injection, variance in parameter update is initially positively correlated with field 1654 amplitude (blue). When a small amount of noise is added, fields with a larger mean amplitude show 1655 a smaller variance in change in parameter while fields with a smaller amplitude show higher vari-1656 ance. Conversely, when the magnitude of noise is further increased (purple), fields with a higher 1657 amplitude show higher variance in its parameters. (C) The correlation between mean amplitude and the magnitude of the readout weights (sum over all actions for squared actor weights and squared 1658 critic weights) is high and positively correlated when the noise magnitude is low. This correla-1659 tion decreases and becomes weakly positive when $\sigma_{noise} = 0.001$. This supports the claim that in 1660 the low noise regime, fields with a high amplitude are more involved in policy learning and hence 1661 drift less or are more stable to maintain performance integrity. (D) Population vector correlation 1662 decreases at a faster rate than the similarity matrix when noise magnitude increases. (E) Represen-1663 tation similarity correlation decreases as the noise magnitude increases, but at a slower rate than PV 1664 correlation. (F) Proportion of fields that are active (average fraction of fields with firing rate less than 1665 0.05, 0.1,0,25) continues to increase with higher noise magnitude. (G) Introducing Gaussian noise 1666 with zero mean and variance N(0, 0.00025) to place field parameters during updates $\theta_{t+1} = \theta_t + \xi_t$ caused each place field's center, firing rate and width to fluctuate as trials progressed. See App. D 1668 for details. This causes each field's spatial selectivity to change over time. Specifically, each field's centroid (λ) shifted from its initialized location, firing rates fluctuated (α^2) causing fields to gain 1669 or lose selectivity, and most fields increased in size (σ^2) while some did not. The first two were 1670 observed by Qin et al. (2023) who analyzed Gonzalez et al. (2019). Each color corresponds to the 1671 dynamics of a specific field, with 5 example fields shown. 1672



1700 Supplementary Figure 9: Noisy place field parameter update replicates drift dynamics seen in 1701 neural data. (A) Place field spatial selectivity changes over days across four mice. Each place fields' centroid positions were sorted according to day 5, 20 and 35. Figure adapted from Fig. 3E-G, 1702 Ziv et al. (2013). (B) Place fields selectivity similarly changes across trials, after stable navigation 1703 performance was attained at trial 25,000. Each place field's centroid position was sorted according 1704 to trial 25,000, 125,000 and 195,000. As trials progress, spatial selectivity becomes distinctively 1705 different similar to Ziv et al. (2013) and Fig. 1G, de Snoo et al. (2023). (C) Probability distributions 1706 of centroid shifts along a 1D track at six (left, adapted from Fig. 3D, Ziv et al. (2013)) and three 1707 (right, adapted from Fig. 5H, Qin et al. (2023) who analyzed Gonzalez et al. (2019) data) different 1708 time intervals. Similar centroid shift away from the initialized position is also observed in Fig. 1709 4B, Geva et al. (2023). (D) When no Gaussian noise is added to place field parameters (α, λ, σ), 1710 place field optimization alone does not cause centroids to shift as in neural data. Instead, adding 1711 small Gaussian noise ($\sigma_{noise} \in \{0.0001, 0.00025, 0.0005\}$) replicates the gradual shift in centroid 1712 position across trials (25,100 to 125,000). When the noise magnitude is high e.g. $\sigma_{noise} \ge 0.001$, 1713 centroids shift rapidly to a new location, similar to the random shuffle or null hypothesis used in Geva et al. (2023); Qin et al. (2023); Ziv et al. (2013). (B-D) Analysis was done for 64 place fields 1714 aggregated over 10 agents initialized with different seeds to have 640 fields in total. (E) Graph 1715 topology of place field activity in a 1D track show clustering of fields according to place encoding 1716 (PC) or end cell (EC, fields found at the end of track). Figure adapted from Fig. 4C, Gonzalez et al. 1717 (2019). (F) Example graph topology for one agent with N = 64 place fields with Gaussian noise 1718 $\sigma_{noise} = 0.00025$ added to field parameters. Each node indicates a place field's centroid position 1719 across learning, and the edge is weighted by the normalized (between 0 to 1) cosine distance between 1720 each node that is less than 0.55. Red, green, blue, orange, black nodes indicate centroids initialized 1721 at the reward, start, end of track near the reward, end of track near the start locations and the middle 1722 of the track respectively. As learning progressed, the cosine distance between each centroid changed 1723 and the ensemble representation rotated. Nevertheless, fields encoding the reward, start, and track were fairly stably as seen in Gonzalez et al. (2019), and the greater separation of clusters support the 1724 phenomenon where a high density of fields emerge at the reward and start locations. 1725 1726

170



Supplementary Figure 10: Influence of field width and number of fields on agent performance. (A) Fields initialized with $\sigma = 0.1$ and (B) $\sigma = 0.05$. Policy learning is slower when initialized with a smaller field width. (C) Influence of field parameter optimization on the average maximum cumulative reward (left) and trial at which agent achieves cumulative discounted reward of 45 and above for the previous 300 trials (right). Correlation plot shows the p-value for a pairwise t-test performed to determine the influence of fields parameters on learning performance.



Supplementary Figure 11: Influence of noise on new target learning performance in 1D track. Increasing the number of place fields (N) and field widths (σ) led to a general increase in new target learning performance. When no noise was injected to field parameters ($\sigma_{noise} = 0.0$, blue), most agents struggled to learn to navigate to new targets and seem to be stuck in a local minima. Instead, noise magnitude of $\sigma_{noise} = 0.0005$ allowed agents to maximize rewards throughout the 250,000 trials. Increasing the noise magnitude beyond this ($\sigma_{noise} = 0.001$) negatively affected the agent's target learning performance, especially when the number of fields were low.



Supplementary Figure 12: Influence of noise on learning performance in 2D arena with 1871 an obstacle. (A) Agents started at the same location $x_{start} = (0.0, 0.75)$ and had to nav-1872 igate to a target that changed to a new location every 50,000 trials following the sequence 1873 $(x_r \in [(0.75, -0.75), (-0.75, 0.75), (0.75, 0.75), (-0.75, -0.75)])$. Increasing the noise magni-1874 tude improved new target learning performance. (B) Agents learned to navigate to a target at $x_r = (0.75, 0.0)$ from a start location $x_{start} = (-0.75, 0.0)$ with an obstacle with coordinates 1875 $(x_{min} = -0.2, x_{max} = 0.2, y_{min} = -1.0, y_{max} = 0.5)$ for the first 50,000 trials. After which, the 1876 location of the obstacle was shifted up to $(x_{min} = -0.2, x_{max} = 0.2, y_{min} = -0.5, y_{max} = 1.0)$ 1877 while the start and target location was the same. Agents with a noise magnitude $\sigma_{noise} = 0.00025$ 1878 showed the highest average reward maximization performance followed by $\sigma_{noise} = 0.0005$. A 1879 high noise magnitude ($\sigma_{noise} = 0.001$) disrupted learning performance while agents without noisy 1880 field updates ($\sigma_{noise} = 0.0$) did not learn to navigate around the new obstacle. Note that field ampli-1881 tudes and widths were clipped to be between $[10^{-5}, 2]$ and $[10^{-5}, 0.5]$ respectively to ensure the Σ 1882 covariance matrix in 2D place fields remained valid for matrix inversion. Performance was averaged 1883 over agents initialized with different number of 2D place fields ($N \in \{64, 144, 256, 576\}$) with the 1884 diagonals of the field width initialized with $\Sigma = 0.01$ and constant amplitude $\alpha = 1.0$, over 30 1885 different seeds. Shaded area is 95% CI.

- 1886
- 1887
- 1000
- 1889



Supplementary Figure 13: Using the same learning rates for the place field parameters and 1928 actor-critic recovers the same phenomena of a high field density emerging at reward location 1929 followed by the start location, and field elongation against the agent's trajectory. (A) Each place 1930 field's amplitude, center and width were sampled from a uniform distribution of [0,1], [-1,1], [1e-1931 5,0.1] respectively to model heterogeneous place field distribution. After learning, a high density 1932 (number) of fields emerged at the start (green dash) and reward (red area) location, similar to Fig. 1933 1B (right) and Sup. Fig. 2B. This phenomenon is consistent across different numbers of place fields. 1934 Shaded area is 95% CI over 50 different seeds. (C) In a 2D arena with obstacles, place fields elongate 1935 from the reward location (red circle) back to the start location (green circle), while narrowing along 1936 the corridor with an obstacle (gray), similar to Fig. 2F. Learning rates for the actor, critic and place 1937 field parameters were $\eta = \eta_{\theta} = 0.0005$.

- 1938
- 1939
- 1940



Supplementary Figure 14: Center-surround place fields reproduces the emergence of a high density of fields at the reward location. (A) Example of 16 center-surround fields uniformly distributed before (left) and after learning for 10,000 trials (right), with the learning rates for the center-surround place field parameters and policy network being the same ($\eta = \eta_{\theta} = 0.001$). Place fields near the reward shifted to the reward location while others elongated from the reward location back to the start location similar to Fig. 2C (bottom row). (B) A high field density (gKDE(COM)) and mean firing rate $(\sum \phi(x))$ emerged at the reward location for N = 16 (left) and N = 64 (right) when using center-surrounds fields. However, we do not see a high density emerging at the start location robustly. Further analysis is needed to verify if the representations learned by Gaussian basis functions and center-surround fields (difference of Gaussians) are similar, and if not why. Shaded area is 95% CI for 10 seeds.



Supplementary Figure 15: Difference in policy convergence when backpropagating temporal difference error through the actor and/or critic weights to optimize place field parameters. We evaluated the speed of policy learning when optimizing place field parameters using (1) the actor weights W^{π} multiplied by the normalized action vector $\tilde{g}_t = q_t - P$ and the critic weights w^{v} (blue) (2) only the the actor weights multiplied by the normalized action vector (orange) (3) only the critic weights (green) (4) direct feedback of the TD error to modulate field parameters instead of backpropagating through the actor or critic weights, making it more biologiclly plausible (red). The combined objective used for place field parameter optimization achieved the fastest policy convergence when the number of fields was low $(N = \{8, 16, 32\})$ (blue). With more fields, using the critic weights (green) was almost as effective as the combined objective. Optimizing place field parameters using only the actor weights led to the slowest policy convergence (orange). Surprisingly, direct feedback of the TD error to modulate place field parameters shows the 2nd fastest policy convergence. Additional analysis is needed to determine the nature of representations learned by all four methods. Shaded area indicates 95%CI over 30 random seeds with place field amplitudes and widths uniformly initialized between [0, 1] and $[10^{-5}, 0.1]$ respectively.