PLACE FIELD REPRESENTATION LEARNING DURING POLICY LEARNING

M Ganesh Kumar^{1,2}, Blake Bordelon^{1,2}, Jacob A. Zavatone-Veth^{2,3}, Cengiz Pehlevan^{1,2,4}

¹John A. Paulson School of Engineering and Applied Sciences
²Center for Brain Science
³Society of Fellows
⁴Kempner Institute for the Study of Natural and Artificial Intelligence
Harvard University
Cambridge, MA 02138
{mganeshkumar@seas,blake_bordelon@g,jzavatoneveth@fas,cpehlevan@seas}.harvard.edu

ABSTRACT

As rodents navigate in a novel environment, a high place field density 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 reinforcement learning, whereby the Temporal Difference (TD) error modulates place field representations to improve policy learning. Place fields are modeled using Gaussian radial basis functions to represent spatial information, and directly synapse to an actor and 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 reward dependent objective. We demonstrate that this framework unifies three disparate phenomena observed in navigation experiments. Furthermore, we show that these place field representations improve policy convergence when learning to navigate to a single target and relearning 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 a hippocampal neuron is maximal and robust across trials (O'Keefe & Dostrovsky, 1971; O'Keefe, 1978). Classically, each neuron has a unique spatial receptive field such that the population activity can describe an animal's allocentric position within the environment (Moser et al., 2015). Ablation 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 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 backward against the trajectory (Mehta et al., 1997; Priestley et al., 2022), and individual field's spatial selectivity continues to change or "drift" even when animals demonstrate stable behavior (Geva et al., 2023; Krishnan & Sheffield, 2023; Kentros et al., 2004; Mankin et al., 2012; Ziv et al., 2013). Although disparate mechanisms have been proposed to model these phenomena, a framework that can unify these phenomena and clarify their computational role remains elusive.

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:

• 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 or "state", which feed into the actor-critic network in the second layer that uses these representations to learn actions that 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: (1) the emergence of high place field density at rewards, (2) elongation of fields against the trajectory, and (3) 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, the mean population 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.

2 RELATED WORKS

Anatomically constrained architecture for navigation. Learning to navigate involves the hippocampus encoding spatial information and its strong glutamatergic projections to the striatum (Lisman & Grace, 2005; Floresco et al., 2001). The ventral and dorsal regions of the striatum are associated with value estimation and stimulus-response associations, functioning similarly to a critic and an actor, respectively (Niv, 2009; Joel et al., 2002; Houk et al., 1994). Additionally, dopamine neurons in the Ventral Tegmental Area influence plasticity in the striatal synapses (Reynolds et al., 2001; Russo & Nestler, 2013). This anatomical insight has led to the design of a biologically plausible navigation model, where place fields connect directly to an actor-critic framework, and synapses are modulated by the TD error (Arleo & Gerstner, 2000; Foster et al., 2000; Frémaux et al., 2013; Brown & Sharp, 1995; Kumar et al., 2022). Recent evidence shows direct dopaminergic projections to the hippocampus to modulate place cell activity, strengthening the case for navigation models with adaptive place fields (Palacios-Filardo & Mellor, 2019; Krishnan et al., 2022; Kempadoo et al., 2016; Sayegh et al., 2024). How upstream information from the entorhinal cortex influences place field representations for policy learning needs clarity (Fiete et al., 2008; Bush et al., 2015). As new experiments challenge the canonical definition that a place cell only has one place field (Eliav et al., 2021), we study spatial representational learning using Gaussian place fields, instead of place cells.

Field density increases near reward locations. Density traditionally refers to the number of field centers of mass in a location. However, we also consider changes in the mean population firing rate, which includes variations in each field's width and amplitude. As animals learn to navigate towards a reward, a high density of place fields emerge at reward locations (Gauthier & Tank, 2018; Lee et al., 2020; Sosa et al., 2023). Reward location based reorganization was observed in hippocampal CA1 and not in CA3 (Dupret et al., 2010). Interestingly, a recent study showed that place fields initially coding for reward shifted backwards against the trajectory causing a decrease in reward coding fields, suggesting of a representation predictively coding for reward (Yaghoubi et al., 2024).

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

Fields drift during stable behavior. After animals reach a certain performance criterion in navigating to a reward location, the spatial selectivity of individual place fields changes across days, even though animals exhibit stable behavior (Kentros et al., 2004; Mankin et al., 2012; Ziv et al., 2013; Geva et al., 2023; de Snoo et al., 2023). A proposal is that these fields continue to drift within a degenerate solution space while the overall representational manifold or the chosen performance metric remains stable (Qin et al., 2023; Pashakhanloo & Koulakov, 2023; Masset et al., 2022; Kappel

et al., 2015; Rokni et al., 2007). Another proposal is that compensatory synaptic plasticity adjusts the readout to maintain stable decoding over time (Rule et al., 2020; Rule & O'Leary, 2022). However, a model that demonstrates stable navigation learning behavior with drifting fields is absent, and the functional role of drift remains unclear.

3 TASK AND MODEL SETUP

Most navigational experiments involve an animal moving from a start location to a target location to receive a reward, either in a one-dimensional (1D) track or a two-dimensional (2D) arena. Similarly, our agents receive their true position at every time step (t) described by the variable (scalar x_t in 1D, vector x_t in 2D), and have to learn a policy (π) that specifies the actions to take (g_t) to move from a start location (e.g. $x_{start} = -0.75$, Fig. 1A green dash) to a target with reward values following a Gaussian distribution ($x_r = 0.5$, $\sigma_r = 0.05$, Fig. 1A red area). The agent outputs a one-hot vector g_t (left or right in 1D and left, right, up or down in 2D), which causes its motion to be discrete, similar to a trajectory in a grid world. To model smooth trajectories in a continuous space as an animal's behavior (Foster et al., 2000; Frémaux et al., 2013; Kumar et al., 2022; 2024b), we use a low-pass filter to smooth g_t using a constant $\alpha_{env} = 0.2$ after scaling for maximum displacement using $v_{max} = 0.1$:

$$x_{t+1} = x_t + \bar{g}_t, \ \bar{g}_{t+1} = (1 - \alpha_{env})\bar{g}_t + \alpha_{env}v_{max}g_t.$$
(1)

To track an agent's reward maximization performance during navigational learning we compute the cumulative discounted reward ($G = \sum_{t=0}^{T} \sum_{k=0}^{T} \gamma^k r_{t+1+k}$) for each trial using $\gamma = 0.9$ as the discount factor, which is similar to tracking the cumulative reward. The trial is terminated when the maximum trial time is reached T_{max} or when the total reward achieved $\sum_{t=0}^{T} r_{t+1}$ reaches a threshold R_{max} .

3.1 PLACE FIELDS AS SPATIAL FEATURES

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

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

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

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:

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

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_{ii}^{π} were initialized by sampling from a normal distribution $\mathcal{N}(0, 10^{-5})$.

3.3 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 standard actor-critic algorithm using the temporal difference residual (refer to App. A):

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

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

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):

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

with $\nabla_{w^v} \mathcal{J}^{TD} = 0$ and $\nabla_{W^{\pi}} \mathcal{L} = 0$. We estimate all parameter gradients online, and provide the explicit update equations for each parameter in App. A. The learning rates for the actor-critic and place field parameters can be the same (Fig. 13). For theoretical analysis, we assume a separation of timescales between learning the actor-critic weights and updating place field parameters (App. B).

4 **Results**

4.1 A HIGH DENSITY OF FIELDS EMERGES NEAR THE REWARD LOCATION

We first examine the neural phenomenon where a high field density emerges at the reward location. Field density is defined by the distribution of field centers of mass (COM) (Lee et al., 2020), which we estimate using Gaussian kernel smoothing. Figure 1C shows how our agent's track occupancy (p(x)), field density (d(x)), mean firing rate (f(x)), and individual field's spatial selectivity $(\phi(x))$ change when learning to navigate in a 1D track from the start $x_{start} = -0.75$ 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 higher proportion of time at the start location with sporadic exploration towards the reward. Despite this behavior, a high field density and mean firing rate emerges at the target from a homogeneous field population. Individual fields at the reward location shift closer to the target (Fig. 1), as seen in Sosa et al. (2023), in contrast to fields at non-rewarded locations. As learning progresses and the agent spends a higher proportion of time at the reward, field density and mean firing rate at the start location also begins to rise slightly, replicating the two-peaked field distribution in Gauthier & Tank (2018), with fields shifting backwards towards the start as in Yaghoubi et al. (2024). A high density at the reward location followed by the start location robustly emerges in heterogeneous place field populations when all the field parameters $(\Delta\lambda, \Delta\alpha, \Delta\sigma)$ are optimized (Fig. 1C right, 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. 1D). When optimizing all the field parameters in a homogeneous population, a high field density rapidly emerges at the reward location to increase goal representation (number of COM within 0.25 unit radius from target center) as seen in (Dupret et al., 2010), followed by gradual reorganization of field density backward along the agent's trajectory.

Interestingly, increasing the number of fields in a heterogeneous place field population reduced the average density and mean firing rate (Fig. 1B, Fig. 1) that emerges near the reward location. This is



Figure 1: Fields shift towards and amplify at reward location. (A) The task is to navigate from the start (green dash) to the target (red area) to receive rewards. The agent has N place fields (blue) which synapse to an actor (red) and critic (green). The TD error δ modulates parameter updates. (B) When initialized with a heterogeneous field population, the enhancement of average field density (d(x)) at the reward location x_r compared to non-reward location x' decreases as the number of fields increases. The density decreases when the reward magnitude (R_{max}) decreases, and reward location's size (R_{size}) increases. (C) Example change in field centers for an agent on a 1D track when only optimizing field centers ($\Delta\lambda$). Initially (T = 1000), the agent spends a high proportion of time $(p_{RM}(x))$ at the start location while a high field density (d(x)) and mean firing rate (f(x))emerges at the reward location. As learning proceeds, the agent spends a higher time at the reward location with field density and mean firing rate increasing at the start location (T = 12000). (Right) A high field density and mean firing rate emerge at the reward and start locations for agents initialized with a heterogeneous population and when all parameters are optimized $(\Delta\lambda, \Delta\alpha, \Delta\sigma)$. (D) Example change in field centers for an agent in a 2D arena with an obstacle (gray). In the early learning phase, field centers (black dots) shift to the target, causing a high density to emerge at the reward (10 agents, right). In the later learning phase, the rest of the centers align along the trajectory. The start, reward locations and radius for goal representation (G.R.) are marked by green, red and blue circles. (E) Example field dynamics when an agent (N = 512) navigates a 1D track. (Left) Fields initialized before ($\lambda_i = 0.5$, blue) and after ($\lambda_i = 0.6$, orange) the target move forward and backward respectively, increasing the density near the target. (Right) Fields closest to the reward $(\lambda_i = 0.5; \text{ green})$ show rapid amplification compared to other fields $(\lambda_i = -0.75, 0.0; \text{ blue and})$ orange). The first order perturbative prediction (theory) provides a good approximation. Shaded area and error bars are 95% CI over 50 seeds.

because as the number of fields increase, the agent goes into a weak feature learning regime (Fig. 4) in which feature learning does not contribute to additional advantage. Conversely, the density and mean firing rate are proportional to the reward magnitude, and inversely proportional to the reward location width as a narrower target might require higher discriminability for the agent to maximize rewards. To understand why place fields exhibit these dynamics, we perform a perturbative approximation to the place field parameter changes under TD learning updates (Menache et al., 2005; Bordelon et al., 2024). 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:

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

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



Figure 2: Fields elongate against the trajectory. (A-B) The Reward Maximization (RM), Successor Representation (SR) and Metric Representation (MR) algorithms cause (A) field sizes to increase and (B) center of mass to shift backwards against the trajectory in a 1D track. Field changes were normalized separately to be between 0 to 1 for visualization. (C) All agents initially spend a high proportion of time at the start location, and later learn to dwell at the target (black). Individual SR fields and mean firing rate (red) closely track the proportion of time the agent spends in a location (top). MR fields reorganize only at the start location (middle). Conversely, individual RM fields and mean firing rate 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 (bottom). (D) SR agents show a consistently high, positive correlation (blue) between mean firing rate and proportion of time spent in a location. MR agents' show a non-monotonic increase in correlation (green). Conversely, the RM agents' mean firing rate and time spent at a location become anti-correlated before becoming positively correlated (orange). (E) The SR and RM mean firing rates (blue) become anti-correlated before becoming positively correlated at the later learning phase, while the SR and MR fields align momentarily before de-correlating (orange), and the RM and MR fields become anti-correlated (green). (F) Example change in field selectivity by SR (top), MR (middle), and RM (bottom) agents in a 2D arena with an obstacle. The RM agent's field elongation is more pronounced than the SR and MR agents. Summary statistics in Fig. 6. Shaded area is 95% CI over 10 seeds.

compared to locations with a low value. In addition to the value of a location, the agent's start location (modeled as a Gaussian with mean $\bar{\mu}_x = -0.75$ and spread σ_x) and the mean field center location $\bar{\lambda}$ over time under the policy influence each field's displacement. As the reward location is visited frequently, we expect $\bar{\lambda} \approx 0.5$. As the term within the square bracket changes sign depending on the field location, only the fields near the reward location will shift towards the reward, while the rest of the fields will move towards the start location. Due to these influences, the field density at the reward location will increase first followed by a gradual increase in start location (Fig. 1C,E). Additional approximations are needed to model the agent's trajectory and improve the simulationtheory fit for place field centers (App. B). A similar perturbative analysis for amplitudes yields

$$\alpha_i(t) - \alpha_i(0) \approx 2\frac{\eta_\alpha}{\eta} w_{v,i}^2(t) , \ \eta_\alpha \ll \eta , \tag{8}$$

where $\eta_{\alpha} = 0.0001$ is the learning rate for the α 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 start and reward location with field amplification at the reward location.

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. This behavior suggests predictive coding for future occupancy, which can be learned through Hebbian association of fields (Mehta et al., 2000), or through the successor representation (SR) algorithm, which minimizes state prediction error for each place field to learn the transition probabilities (Stachenfeld et al., 2017). Here we show that both our RM agent and a reward independent Metric Representation (MR) agent recapitulate field elongation in a 1D track. For comparison, we developed two agents: A) SR agent that learns the transition probabilities in parallel to policy learning (Fig. 5A). The SR agent has a similar architecture to our (RM) agent (Fig. 1A), with two key differences: 1) It has one set of place fields with fixed parameters, and only 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 synapses U which are optimized using the SR algorithm (App. C). We compare the learned successor place fields to the learned place fields in our RM model, referring to them henceforth as place fields. B) a Metric Representation (MR) agent (Fig. 5B) that estimates its current coordinates in an environment (z_t) . This representation enables navigation to recalled targets by vector subtraction (Foster et al., 2000; Kumar et al., 2024b). The coordinate readout weights and place field parameters are updated by gradient descent to minimize the path integration-derived TD error $\mathcal{L}^{MR} = \mathbb{E}[\sum_{i=1}^{T} \frac{1}{2}(z_{t+1} - (z_t + a_t))^2]$ (App. D), while only the actor and critic readout weights are updated to learn a policy. This objective optimizes place fields $\phi(x)$ even in the absence of rewards.

All three agents (SR, MR, RM) recapitulate the phenomena seen in (Mehta et al., 1997): on average, place fields increase in size over learning (Fig. 2A), and the center of mass (COM) shifts backwards from their initialized positions (Fig. 2B). However, the place field dynamics evolve differently. All agents initially spend a high proportion of time at the start location and gradually learn a policy to spend a higher proportion of time at the reward (Fig. 2C, Fig. 5E). The SR, by design, tracks the transition probabilities of the agent's policy. Consequently, the SR population mean firing rate $f_{\psi}(x)$ closely aligns with the agent's probability of being in a location p_{SR} , showing a high positive correlation (Fig. 2D, blue). Since the MR representation (Fig. 6D) is only modulated by the agent's displacement a_t , fields reorganize more at the start location since displacement is nonzero, causing a higher mean firing rate. Conversely, displacement becomes zero at the reward location as the agent comes to a stop to maximize rewards, causing low field reorganization at the reward. Hence, MR fields $f_{\phi_{MR}}(x)$ become positively correlated with p_{MR} at the start location, but do not fully align with the agent's time spent at the reward location (Fig. 2D, green). Conversely, during early learning, the RM agent exhibits a high population mean firing rate $f_{\phi_{RM}}(x)$ at the reward location, which contrasts sharply with the proportion of time spent at that location, leading to a highly negative correlation between $f_{\phi_{RM}}(x)$ and p_{RM} (Fig. 2D, orange). Interestingly, in the later phase of learning, $f_{\phi_{RM}}(x)$ and p_{RM} become positively correlated.

The mean firing rates learned by the SR and RM agents become negatively correlated during the early learning phase but become positively correlated at the later learning phase (Fig. 2E, blue). Conversely, the mean firing rate correlation decreases monotonically towards zero for the MR and RM agents (Fig. 2E, green), while the correlation between SR and MR increases due to the alignment at the start location in the early learning phase before becoming uncorrelated in the later learning phase. A similar change in correlation is observed when comparing the individual field selectivity, and the spatial representation similarity matrix (Fig. 5F,G). Hence, while all three algorithms demonstrate similar neural phenomenon, the dynamics of learning these representations are different, with SR and RM agents eventually learn similar spatial representations. In a 2D arena with an obstacle, the three agents show field elongation against the movement direction (Fig. 2, 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 and field elongation by MR is more pronounced at the start location.

4.3 STABLE NAVIGATION BEHAVIOR WITH DRIFTING FIELDS

The third phenomena that the model captures has been described as representational drift, where the agent demonstrates stable behavior but the spatial selectivity of individual place fields changes over time (Fig. 3A, Fig. 8G), as seen in Ziv et al. (2013). Although our agent uses a stochastic policy, both the navigation behavior (Fig. 3E, blue) and the population vector (PV) correlation (Fig. 3C, blue) are extremely stable. To drive larger variability in the representation, we introduced Gaussian noise to the field parameter updates at every time step (App. E). Increasing the noise magnitude led to a faster decrease in PV correlation but also disrupted agents' policy convergence for magnitudes greater than 10^{-3} (Fig. 3E, 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 while the PV correlation decreases rapidly (Fig. 3C,E). This demonstrates that agents can optimize their policies to maintain stable behavior even though individual spatial selectivity is



Figure 3: Stable representation similarity and anchor fields facilitate consistent behavior. (A) Injecting Gaussian noise with magnitude $\sigma_{noise} = 0.0001$ into field parameters causes individual field's spatial selectivity to change across trials while (B) the representation similarity matrix (dot product of population activity) remains stable. (C) Injecting higher noise magnitudes $(\sigma_{noise} = 0.0, 0.0001, 0.001)$ leads to a faster decrease in population vector correlation (R_{PV}) across trials while (D) the similarity matrix correlation (R_{RS}) decreases at a slower rate. (E) Agents' reward maximization performance (G) remains fairly stable when the noise magnitude increases. Beyond $\sigma_{noise} = 0.001$, performance becomes highly unstable. Black dash indicates the trial at which PV and similarity matrix correlation was measured from. (F) Normalized variance in field parameters ($\theta = {\alpha, \lambda, \sigma}$) 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 amplitude experiences a greater change in its parameters. When $\sigma_{noise} = 0.0001$, we see the opposite trend, where fields with a larger amplitude are more stable than fields with a smaller amplitude. Refer to Fig. 8 for other σ_{noise} values. Shaded area is 95% CI over 10 seeds.

changing. Interestingly, the spatial representation similarity matrix remains more stable than PV correlation (Fig. 3B,D), even with a higher noise magnitude (Fig. 3D), 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 reward maximization behavior, representation similarity and PV correlation to change at similar rates (Fig. 7), which is not as consistent with experiments (See Fig. 9 for comparisons to data).

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

4.4 PLACE FIELD REORGANIZATION IMPROVES POLICY CONVERGENCE

As the reward-maximizing model recapitulates experimentally-observed changes in place fields, it 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 ablations are particularly important due to the parameter degeneracies in the model: one can trade off the place field amplitudes and the critic and actor weights. We first considered the task of navigating to a single fixed target. Agents with fixed place fields attained the lowest navigational performance with cumulative reward G plateauing at G = 33 (Fig. 4A), and showed the slowest policy convergence even as the number of fields increased (Fig. 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 field amplitudes (α) contributed to



Figure 4: Field reorganization and noisy updates improve target learning. (A) Optimizing all three field parameters, amplitude, width and center of randomly distributed fields allowed agents $(N = 16, \sigma = 0.1)$ to attain the highest cumulative discounted reward (G), while fields with fixed field parameters attained the lowest. (B) Optimizing place field widths (σ) , followed by field amplitudes (α) and lastly field centers (λ) caused the biggest decrease in the number of trials needed for policy convergence $(T_{G>45}, \text{ attain a running average of } G = 45 \text{ over } 300 \text{ trials})$. As the number of fields increased, the number of trials needed for policy convergence decreased and the computational advantage afforded by field optimization extinguished. (C) Agents need to navigate to a target that changed after 50,000 trials $x_r = \{0.5, 0.0, 0.75, -0.25, 0.5\}$. Without noisy field parameter updates, 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 reward for subsequent targets, with $\sigma_{noise} = 0.0005$ (red) demonstrating the highest improvement. Shaded area is 95% CI over 50 seeds.

the next most significant improvement (Fig. 4A-B). Interestingly, place field center (λ) optimization did not contribute to a significant improvement in performance, and in fact caused a decrease in reward maximization performance and speed of policy convergence when optimized together with the amplitude parameter. Hence, optimizing field widths followed by amplitudes and lastly centers significantly improved agent's reward maximization performance and increased the speed of policy convergence. Optimizing field parameters using the auxiliary metric representation objective, inspired by Fang & Stachenfeld (2023), marginally improved policy learning (Fig. 15). However, as the number of place fields increase (Fig. 4B), the computational advantage afforded by place field optimization extinguishes. Nevertheless, optimizing all the parameters in a small number of fields, e.g. 8, leads to a similar rate of policy convergence than with a larger number of randomly initialized fields e.g. 128, which hints that representation flexibility could allow efficient learning in systems with few neurons.

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

5 DISCUSSION

We present a two-layer navigation model which uses tunable place fields as feature inputs to an actor and a critic for policy learning. The parameters of the place fields and the policy and value function learn to maximize rewards using the temporal difference (TD) error. Our simple reinforcement learning model reproduces three experimentally-observed neural phenomena: (1) the emergence of 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 error, number of place fields and noise magnitudes influenced place field representations. Lastly, we demonstrate that learning place field representations with noisy field parameters improves the rate of policy convergence when learning single and multiple targets.

The proposed reinforcement learning model might be amenable to theoretical analysis (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 account for field elongation against the trajectory, which can be contrasted with the successor representation algorithm (Raju et al., 2024; Kumar et al., 2024b). As field dynamics are different in these two models, they could be distinguished by experiments that track fields over the full course of learning (Fig. 2C-E, Fig. 6). Furthermore, place field width and amplitude optimization increases maximum cumulative reward and accelerates policy convergence (Fig. 4A-B).

Most models that characterized representational drift were not studied under the context of navigational policy learning (Masset et al., 2022; Pashakhanloo & Koulakov, 2023; Ratzon et al., 2024). We showed that increasing the noise magnitudes caused different drift regimes (Fig. 3F; Fig. 9D), and at very high noise levels navigation behavior started to collapse (Fig. 3C, Fig. 7). Importantly, we showed that fields in the noisy regime allowed agents to consistently learn new targets in both 1D (Fig. 4C) and 2D (Fig. 12A-B) environments, without getting stuck in local minima. The biological origins of adding noise to place field parameters can be attributed to noisy synaptic plasticity mechanisms (Mongillo et al., 2017; Kappel et al., 2015; Attardo et al., 2015). Other mechanisms such as unstable dynamics in downstream networks (Sorscher et al., 2023) and modulatory mechanisms such dopamine fluctuations (Krishnan & Sheffield, 2023) could adaptively control drift rates. A difficult experiment that could directly verify our model is to induce or constrain place 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 rates needs to be explored. The current model provides a starting point for this investigation.

5.1 LIMITATIONS AND FUTURE WORK

The proposed model is not without limitations. First, we modeled single peaked place fields instead of the complex representations resulting from single "place" cells, which can be multi-field and multi-scale. Nevertheless, the proposed online reinforcement learning framework is general enough to accommodate other models of place cell description (Mainali et al., 2024; Sorscher et al., 2023) e.g. Fig. 14, and can be extended to study representation learning in other brain regions e.g. medial entorhinal (Boccara et al., 2019; Wen et al., 2024) or posterior parietal (Suhaimi et al., 2022) cortex. Next, place field parameters are optimized by backpropagating the temporal difference error through the actor and critic components (Fig. 15). Since the motivation was to develop a normative model whose objective was to maximize rewards, this was a reasonable starting point. However, this model must be extended using biologically-plausible learning rules (Miconi, 2017; Murray, 2019; Lillicrap et al., 2016; Nøkland, 2016) before it can in any way be considered mechanistic (Lee et al., 2024; Starkweather & Uchida, 2021; Krishnan et al., 2022; Kempadoo et al., 2016; Edelmann & Lessmann, 2018). Although we explored a simple non-reward-dependent objective to drive place field reorganization, extending the model to other auxiliary objectives (Low et al., 2018; Schaeffer et al., 2022) to understand their influence in representation learning for policy learning is the next step. While our computational experiments successfully demonstrated the model's effectiveness in reproducing three disparate phenomena, further work should test its robustness across other reinforcement learning algorithms e.g. policy gradient (Kumar & Pehlevan, 2024). Additionally, we need to explore how place field reorganization scales in larger, more complex environments (Hill et al., 2020; Lin et al., 2023; Nieh et al., 2021; Kumar et al., 2024b) beyond the few environments we considered. Lastly, we need to quantitatively compare the representation alignment (Lampinen et al., 2024; Cloos et al., 2024) between our model's place field dynamics and experimental data.

CODE AVAILABILITY

The code for our agents and to reproduce all figures in this paper is available at:

https://github.com/Pehlevan-Group/placefield_reorg_agent

AUTHOR CONTRIBUTIONS

MGK and CP conceptualized and designed the study. MGK and BB performed the theoretical analysis. MGK performed the simulations and wrote the original draft. BB, JZV and CP revised the manuscript.

ACKNOWLEDGMENTS

We would like to thank Albert Lee, Lucas Janson, Farhad Pashakhanloo, Shahriar Talebi, Paul Masset, as well as the members of the Pehlevan, Ba, Janson and Murthy labs for useful insights. We also appreciate the discussions during the Analytical Connectionism Summer School 2024. This research was supported in part by grants NSF PHY-1748958 and PHY-2309135 to the Kavli Institute for Theoretical Physics (KITP). MGK and CP are supported by NSF Award DMS-2134157. CP is further supported by NSF CAREER Award IIS-2239780, and a Sloan Research Fellowship. BB is supported by a Google PhD Fellowship. JAZV is supported by a Junior Fellowship from the Harvard Society of Fellows. This work has been made possible in part by a gift from the Chan Zuckerberg Initiative Foundation to establish the Kempner Institute for the Study of Natural and Artificial Intelligence.

REFERENCES

- 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. *Nature neuroscience*, 25(8):1082–1092, 2022.
- Angelo Arleo and Wulfram Gerstner. Spatial cognition and neuro-mimetic navigation: a model of hippocampal place cell activity. *Biological cybernetics*, 83(3):287–299, 2000.
- Alessio Attardo, James E Fitzgerald, and Mark J Schnitzer. Impermanence of dendritic spines in live adult cal hippocampus. *Nature*, 523(7562):592–596, 2015.
- Charlotte N Boccara, Michele Nardin, Federico Stella, Joseph O'Neill, and Jozsef Csicsvari. The entorhinal cognitive map is attracted to goals. *Science*, 363(6434):1443–1447, 2019.
- Blake Bordelon, Paul Masset, Henry Kuo, and Cengiz Pehlevan. Loss dynamics of temporal difference reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- 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. *Hippocampus*, 5 (3):171–188, 1995.
- Daniel Bush, Caswell Barry, Daniel Manson, and Neil Burgess. Using grid cells for navigation. *Neuron*, 87(3):507–520, 2015.
- Nathan Cloos, Moufan Li, Markus Siegel, Scott L. Brincat, Earl K. Miller, Guangyu Robert Yang, and Christopher J. Cueva. Differentiable optimization of similarity scores between models and brains, 2024. URL https://arxiv.org/abs/2407.07059.
- Peter Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural computation*, 5(4):613–624, 1993.
- Peter Dayan and Laurence F Abbott. *Theoretical neuroscience: computational and mathematical modeling of neural systems*. MIT press, 2005.
- Mitchell L de Snoo, Adam MP Miller, Adam I Ramsaran, Sheena A Josselyn, and Paul W Frankland. Exercise accelerates place cell representational drift. *Current Biology*, 33(3):R96–R97, 2023.

- Shibhansh Dohare, J Fernando Hernandez-Garcia, Qingfeng Lan, Parash Rahman, A Rupam Mahmood, and Richard S Sutton. Loss of plasticity in deep continual learning. *Nature*, 632(8026): 768–774, 2024.
- Can Dong, Antoine D Madar, and Mark EJ Sheffield. Distinct place cell dynamics in ca1 and ca3 encode experience in new environments. *Nature communications*, 12(1):2977, 2021.
- David Dupret, Joseph O'neill, Barty Pleydell-Bouverie, and Jozsef Csicsvari. The reorganization and reactivation of hippocampal maps predict spatial memory performance. *Nature neuroscience*, 13(8):995–1002, 2010.
- Elke Edelmann and Volkmar Lessmann. Dopaminergic innervation and modulation of hippocampal networks. *Cell and tissue research*, 373:711–727, 2018.
- 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. *Science*, 372(6545):eabg4020, 2021.
- Ching Fang and Kimberly L Stachenfeld. Predictive auxiliary objectives in deep rl mimic learning in the brain. *arXiv preprint arXiv:2310.06089*, 2023.
- Ila R Fiete, Yoram Burak, and Ted Brookings. What grid cells convey about rat location. *Journal of Neuroscience*, 28(27):6858–6871, 2008.
- Stan B Floresco, Christopher L Todd, and Anthony A Grace. Glutamatergic afferents from the hippocampus to the nucleus accumbens regulate activity of ventral tegmental area dopamine neurons. *Journal of Neuroscience*, 21(13):4915–4922, 2001.
- David J Foster, Richard GM Morris, and Peter Dayan. A model of hippocampally dependent navigation, using the temporal difference learning rule. *Hippocampus*, 10(1):1–16, 2000.
- Loren M Frank, Garrett B Stanley, and Emery N Brown. Hippocampal plasticity across multiple days of exposure to novel environments. *Journal of Neuroscience*, 24(35):7681–7689, 2004.
- Nicolas Frémaux, Henning Sprekeler, and Wulfram Gerstner. Reinforcement learning using a continuous time actor-critic framework with spiking neurons. *PLoS computational biology*, 9(4): e1003024, 2013.
- Matthew PH Gardner, Geoffrey Schoenbaum, and Samuel J Gershman. Rethinking dopamine as generalized prediction error. *Proceedings of the Royal Society B*, 285(1891):20181645, 2018.
- Jeffrey L Gauthier and David W Tank. A dedicated population for reward coding in the hippocampus. *Neuron*, 99(1):179–193, 2018.
- Samuel J Gershman. The successor representation: its computational logic and neural substrates. *Journal of Neuroscience*, 38(33):7193–7200, 2018.
- Samuel J Gershman and Naoshige Uchida. Believing in dopamine. *Nature Reviews Neuroscience*, 20(11):703–714, 2019.
- Nitzan Geva, Daniel Deitch, Alon Rubin, and Yaniv Ziv. Time and experience differentially affect distinct aspects of hippocampal representational drift. *Neuron*, 111(15):2357–2366, 2023.
- Walter G Gonzalez, Hanwen Zhang, Anna Harutyunyan, and Carlos Lois. Persistence of neuronal representations through time and damage in the hippocampus. *Science*, 365(6455):821–825, 2019.
- Felix Hill, Olivier Tieleman, Tamara Von Glehn, Nathaniel Wong, Hamza Merzic, and Stephen Clark. Grounded language learning fast and slow. *arXiv preprint arXiv:2009.01719*, 2020.
- James C. Houk, James L. Adams, and Andrew G. Barto. A Model of How the Basal Ganglia Generate 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/4708.003.0020. URL https://doi.org/10.7551/mitpress/4708.003.0020.

- Daphna Joel, Yael Niv, and Eytan Ruppin. Actor–critic models of the basal ganglia: New anatomical and computational perspectives. *Neural networks*, 15(4-6):535–547, 2002.
- David Kappel, Stefan Habenschuss, Robert Legenstein, and Wolfgang Maass. Network plasticity as bayesian inference. *PLoS computational biology*, 11(11):e1004485, 2015.
- Kimberly A Kempadoo, Eugene V Mosharov, Se Joon Choi, David Sulzer, and Eric R Kandel. Dopamine release from the locus coeruleus to the dorsal hippocampus promotes spatial learning and memory. *Proceedings of the National Academy of Sciences*, 113(51):14835–14840, 2016.
- Clifford G Kentros, Naveen T Agnihotri, Samantha Streater, Robert D Hawkins, and Eric R Kandel. Increased attention to spatial context increases both place field stability and spatial memory. *Neuron*, 42(2):283–295, 2004.
- Seetha Krishnan and Mark EJ Sheffield. Reward expectation reduces representational drift in the hippocampus. *bioRxiv*, 2023.
- Seetha Krishnan, Chad Heer, Chery Cherian, and Mark EJ Sheffield. Reward expectation extinction restructures and degrades cal spatial maps through loss of a dopaminergic reward proximity signal. *Nature communications*, 13(1):6662, 2022.
- M Ganesh Kumar and Cengiz Pehlevan. Place fields organize along goal trajectory with reinforcement learning. *Cognitive Computational Neuroscience*, 2024.
- M Ganesh Kumar, Cheston Tan, Camilo Libedinsky, Shih-Cheng Yen, and Andrew YY Tan. A nonlinear hidden layer enables actor–critic agents to learn multiple paired association navigation. *Cerebral Cortex*, 32(18):3917–3936, 2022.
- M Ganesh Kumar, Shamini Ayyadhury, and Elavazhagan Murugan. Trends innovations challenges in employing interdisciplinary approaches to biomedical sciences. In *Translational Research in Biomedical Sciences: Recent Progress and Future Prospects*, pp. 287–308. Springer, 2024a.
- M Ganesh Kumar, Cheston Tan, Camilo Libedinsky, Shih-Cheng Yen, and Andrew Yong-Yi Tan. One-shot learning of paired association navigation with biologically plausible schemas, 2024b. URL https://arxiv.org/abs/2106.03580.
- Andrew Kyle Lampinen, Stephanie C. Y. Chan, and Katherine Hermann. Learned feature representations are biased by complexity, learning order, position, and more, 2024. URL https: //arxiv.org/abs/2405.05847.
- Jae Sung Lee, John J Briguglio, Jeremy D Cohen, Sandro Romani, and Albert K Lee. The statistical structure of the hippocampal code for space as a function of time, context, and value. *Cell*, 183 (3):620–635, 2020.
- 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): 1574–1586, 2024.
- Timothy P Lillicrap, Daniel Cownden, Douglas B Tweed, and Colin J Akerman. Random synaptic feedback weights support error backpropagation for deep learning. *Nature communications*, 7(1): 13276, 2016.
- Zijun Lin, Haidi Azaman, M Ganesh Kumar, and Cheston Tan. Compositional learning of visuallygrounded concepts using reinforcement. *arXiv preprint arXiv:2309.04504*, 2023.
- John E Lisman and Anthony A Grace. The hippocampal-vta loop: controlling the entry of information into long-term memory. *Neuron*, 46(5):703–713, 2005.
- Ryan J Low, Sam Lewallen, Dmitriy Aronov, Rhino Nevers, and David W Tank. Probing variability in a cognitive map using manifold inference from neural dynamics. *BioRxiv*, pp. 418939, 2018.
- Nischal Mainali, Rava Azeredo da Silveira, and Yoram Burak. Universal statistics of hippocampal place fields across species and dimensionalities. *bioRxiv*, pp. 2024–06, 2024.

- Emily A Mankin, Fraser T Sparks, Begum Slayyeh, Robert J Sutherland, Stefan Leutgeb, and Jill K Leutgeb. Neuronal code for extended time in the hippocampus. *Proceedings of the National Academy of Sciences*, 109(47):19462–19467, 2012.
- Paul Masset, Shanshan Qin, and Jacob A Zavatone-Veth. Drifting neuronal representations: Bug or feature? *Biological cybernetics*, 116(3):253–266, 2022.
- Mayank R Mehta, Carol A Barnes, and Bruce L McNaughton. Experience-dependent, asymmetric expansion of hippocampal place fields. *Proceedings of the National Academy of Sciences*, 94(16): 8918–8921, 1997.
- Mayank R Mehta, Michael C Quirk, and Matthew A Wilson. Experience-dependent asymmetric shape of hippocampal receptive fields. *Neuron*, 25(3):707–715, 2000.
- Ishai Menache, Shie Mannor, and Nahum Shimkin. Basis function adaptation in temporal difference reinforcement learning. *Annals of Operations Research*, 134(1):215–238, 2005.
- Thomas Miconi. Biologically plausible learning in recurrent neural networks reproduces neural dynamics observed during cognitive tasks. *Elife*, 6:e20899, 2017.
- 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 International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 1928–1937, New York, New York, USA, 20–22 Jun 2016. PMLR. URL https: //proceedings.mlr.press/v48/mniha16.html.
- Gianluigi Mongillo, Simon Rumpel, and Yonatan Loewenstein. Intrinsic volatility of synaptic connections—a challenge to the synaptic trace theory of memory. *Current opinion in neurobiology*, 46:7–13, 2017.
- P Read Montague, Peter Dayan, and Terrence J Sejnowski. A framework for mesencephalic dopamine systems based on predictive hebbian learning. *Journal of neuroscience*, 16(5):1936– 1947, 1996.
- Richard GM Morris, Paul Garrud, JNP al Rawlins, and John O'Keefe. Place navigation impaired in rats with hippocampal lesions. *Nature*, 297(5868):681–683, 1982.
- May-Britt Moser, David C Rowland, and Edvard I Moser. Place cells, grid cells, and memory. *Cold Spring Harbor perspectives in biology*, 7(2):a021808, 2015.
- James M Murray. Local online learning in recurrent networks with random feedback. *Elife*, 8: e43299, 2019.
- Edward H Nieh, Manuel Schottdorf, Nicolas W Freeman, Ryan J Low, Sam Lewallen, Sue Ann Koay, Lucas Pinto, Jeffrey L Gauthier, Carlos D Brody, and David W Tank. Geometry of abstract learned knowledge in the hippocampus. *Nature*, 595(7865):80–84, 2021.
- Yael Niv. Reinforcement learning in the brain. *Journal of Mathematical Psychology*, 53(3):139–154, 2009.
- Arild Nøkland. Direct feedback alignment provides learning in deep neural networks. Advances in neural information processing systems, 29, 2016.
- J O'Keefe. The hippocampus as a cognitive map, 1978.
- John O'Keefe and Jonathan Dostrovsky. The hippocampus as a spatial map: preliminary evidence from unit activity in the freely-moving rat. *Brain research*, 1971.
- Mark G Packard and James L McGaugh. Inactivation of hippocampus or caudate nucleus with lidocaine differentially affects expression of place and response learning. *Neurobiology of learning and memory*, 65(1):65–72, 1996.
- Jon Palacios-Filardo and Jack R Mellor. Neuromodulation of hippocampal long-term synaptic plasticity. Current opinion in neurobiology, 54:37–43, 2019.

- Farhad Pashakhanloo and Alexei Koulakov. Stochastic gradient descent-induced drift of representation in a two-layer neural network. In *International Conference on Machine Learning*, pp. 27401–27419. PMLR, 2023.
- James B Priestley, John C Bowler, Sebi V Rolotti, Stefano Fusi, and Attila Losonczy. Signatures of rapid plasticity in hippocampal ca1 representations during novel experiences. *Neuron*, 110(12): 1978–1992, 2022.
- 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. *Nature Neuroscience*, 26(2):339–349, 2023.
- 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. *Science Advances*, 10(31):eadm8470, 2024.
- Aviv Ratzon, Dori Derdikman, and Omri Barak. Representational drift as a result of implicit regularization. *Elife*, 12:RP90069, 2024.
- John NJ Reynolds, Brian I Hyland, and Jeffery R Wickens. A cellular mechanism of reward-related learning. *Nature*, 413(6851):67–70, 2001.
- Uri Rokni, Andrew G Richardson, Emilio Bizzi, and H Sebastian Seung. Motor learning with unstable neural representations. *Neuron*, 54(4):653–666, 2007.
- Michael E Rule and Timothy O'Leary. Self-healing codes: How stable neural populations can track continually reconfiguring neural representations. *Proceedings of the National Academy of Sciences*, 119(7):e2106692119, 2022.
- Michael E Rule, Adrianna R Loback, Dhruva V Raman, Laura N Driscoll, Christopher D Harvey, and Timothy O'Leary. Stable task information from an unstable neural population. *elife*, 9: e51121, 2020.
- Scott J Russo and Eric J Nestler. The brain reward circuitry in mood disorders. *Nature reviews neuroscience*, 14(9):609–625, 2013.
- 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 hippocampus trigger long-term potentiation and contextual learning. *Nature Communications*, 15 (1):4100, 2024.
- Rylan Schaeffer, Mikail Khona, and Ila Fiete. No free lunch from deep learning in neuroscience: A case study through models of the entorhinal-hippocampal circuit. Advances in neural information processing systems, 35:16052–16067, 2022.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. Highdimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Wolfram Schultz, Peter Dayan, and P Read Montague. A neural substrate of prediction and reward. Science, 275(5306):1593–1599, 1997.
- Ben Sorscher, Gabriel C Mel, Samuel A Ocko, Lisa M Giocomo, and Surya Ganguli. A unified theory for the computational and mechanistic origins of grid cells. *Neuron*, 111(1):121–137, 2023.
- Marielena Sosa, Mark H Plitt, and Lisa M Giocomo. Hippocampal sequences span experience relative to rewards. *bioRxiv*, 2023.
- Kimberly L Stachenfeld, Matthew M Botvinick, and Samuel J Gershman. The hippocampus as a predictive map. *Nature neuroscience*, 20(11):1643–1653, 2017.

- Clara Kwon Starkweather and Naoshige Uchida. Dopamine signals as temporal difference errors: recent advances. *Current Opinion in Neurobiology*, 67:95–105, 2021.
- RJ Steele and RGM Morris. Delay-dependent impairment of a matching-to-place task with chronic and intrahippocampal infusion of the nmda-antagonist d-ap5. *Hippocampus*, 9(2):118–136, 1999.
- Christoph Stöckl, Yukun Yang, and Wolfgang Maass. Local prediction-learning in high-dimensional spaces enables neural networks to plan. *Nature Communications*, 15(1):2344, 2024.
- Ahmad Suhaimi, Amos WH Lim, Xin Wei Chia, Chunyue Li, and Hiroshi Makino. Representation learning in the artificial and biological neural networks underlying sensorimotor integration. *Science Advances*, 8(22):eabn0984, 2022.
- Richard S. Sutton and Andrew G. Barto. Reinforcement learning: An introduction. A Bradford Book, 2018.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- Edward C Tolman. Cognitive maps in rats and men. Psychological review, 55(4):189, 1948.
- Dorothy Tse, Rosamund F Langston, Masaki Kakeyama, Ingrid Bethus, Patrick A Spooner, Emma R Wood, Menno P Witter, and Richard GM Morris. Schemas and memory consolidation. *Science*, 316(5821):76–82, 2007.
- Jane X Wang, Zeb Kurth-Nelson, Dharshan Kumaran, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Demis Hassabis, and Matthew Botvinick. Prefrontal cortex as a meta-reinforcement learning system. *Nature neuroscience*, 21(6):860–868, 2018.
- John H Wen, Ben Sorscher, Emily A Aery Jones, Surya Ganguli, and Lisa M Giocomo. One-shot entorhinal maps enable flexible navigation in novel environments. *Nature*, pp. 1–8, 2024.
- Mohammad Yaghoubi, Andres Nieto-Posadas, Coralie-Anne Mosser, Thomas Gisiger, Émmanuel Wilson, Sylvain Williams, and Mark P Brandon. Predictive coding of reward in the hippocampus. *bioRxiv*, pp. 2024–09, 2024.
- Yaniv Ziv, Laurie D Burns, Eric D Cocker, Elizabeth O Hamel, Kunal K Ghosh, Lacey J Kitch, Abbas El Gamal, and Mark J Schnitzer. Long-term dynamics of ca1 hippocampal place codes. *Nature neuroscience*, 16(3):264–266, 2013.

A DETAILS OF THE PLACE FIELD-BASED NAVIGATION MODEL

The code for initializing and training the model in 1D and 2D environments, along with the code for analyzing neural phenomena and generating all figures, will be available on GitHub upon acceptance.

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

$$\phi_i(x_t) = \alpha_i^2 \exp\left(-\frac{||x_t - \lambda_i||_2^2}{2\sigma_i^2}\right), \qquad (9)$$

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

$$\phi_i(x_t) = \alpha_i^2 \exp\left[-\frac{1}{2}(x_t - \lambda_i)^\top \Sigma_i^{-1}(x_t - \lambda_i)\right], \qquad (10)$$

where Σ_i is a 2x2 covariance matrix, adapted 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

$$\mathcal{J}^G = \mathbb{E}_{\tau \sim \phi_{\theta}, \pi_W \pi} \left[\sum_{t=0}^T \sum_{k=0}^T \gamma^k r_{t+1+k} \right] = \mathbb{E} \left[\sum_{t=0}^T G_t \right], \tag{11}$$

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,

$$\theta_{new} = \theta_{old} + \eta_{\theta} \nabla_{\theta} \mathcal{J}^G \quad , \quad W_{new}^{\pi} = W_{old}^{\pi} + \eta \nabla_{W^{\pi}} \mathcal{J}^G \,, \tag{12}$$

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

$$\nabla_{\theta, W^{\pi}} \mathcal{J}^{G} = \nabla_{\theta, W^{\pi}} \mathbb{E}\left[G(\tau)\right]$$
(13)

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

$$= \int p(\tau|\theta, W^{\pi}) \nabla_{\theta, W^{\pi}} \log p(\tau|\theta, W^{\pi}) G(\tau)$$
(15)

$$= \mathbb{E}\left[\nabla_{\theta, W^{\pi}} \log p(\tau | \theta, W^{\pi}) G(\tau)\right], \qquad (16)$$

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,

$$p(\tau|\theta, W^{\pi}) = p(x_0) \prod_{t=0}^{T} p(x_{t+1}|x_t) \pi(g_t|x_t; \theta, W^{\pi})$$
(17)

$$\log p(\tau|\theta, W^{\pi}) = \log p(x_0) + \sum_{t=0}^{T} \left(\log p(x_{t+1}|x_t) + \log \pi(g_t|x_t; \theta, W^{\pi})\right)$$
(18)

$$\nabla_{\theta, W^{\pi}} \log p(\tau | \theta, W^{\pi}) = \sum_{t=0}^{T} \log \pi(g_t | x_t; \theta, W^{\pi}).$$
(19)

Since the gradients are not dependent on the state transitions, the last line excludes them. Substituting Eq. 19 into Eq. 16 yields

$$\nabla_{\theta, W^{\pi}} \mathcal{J}^{G} = \mathbb{E} \left[\sum_{t=0}^{T} \nabla_{\theta, W^{\pi}} \log \pi(g_{t} | x_{t}; \theta, W^{\pi}) \cdot G_{t} \right],$$
(20)

which completes the full derivation of the policy gradient theorem (Sutton et al., 1999; Sutton & Barto, 2018). 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

$$V_t = \mathbb{E}[G_t], \tag{21}$$

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

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

$$= \mathbb{E}\left[\sum_{t=0}^{T} r_{t+1} + \gamma \sum_{k=0}^{T} \gamma^{k} r_{t+2+k}\right], \qquad (23)$$

$$\mathcal{J}^G = \mathbb{E}\left[\sum_{t=0}^T r_{t+1} + \gamma G_{t+1}\right] = \mathbb{E}\left[\sum_{t=0}^T r_{t+1} + \gamma V_{t+1}\right].$$
(24)

which yields the following self-consistency equation

$$r_{t+1} + \gamma V_{t+1} - V_t \equiv 0, \qquad (25)$$

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

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

$$\mathcal{J}^{A} = \mathbb{E}\left[G_{t} - V_{t}\right] = \mathbb{E}\left[A_{t}\right] = \mathbb{E}\left[\sum_{t=0}^{T} r_{t+1} + \gamma V_{t+1} - V_{t}\right], \qquad (26)$$

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. 11.

$$\nabla_{\theta, W^{\pi}} \mathcal{J}^{A} = \mathbb{E} \left[\sum_{t=0}^{T} \nabla_{\theta} \log \pi(g_{t} | x_{t}; \theta, W^{\pi}) \cdot A_{t} \right].$$
(27)

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

$$\delta_t = r_{t+1} + \gamma v_{t+1} - v_t \,. \tag{28}$$

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

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

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

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 $v(x_{t+1})$

$$\mathcal{L}^{TD} = \mathbb{E}\left[\sum_{t=0}^{T} \frac{1}{2} \left(r_{t+1} + \gamma v(x_{t+1}) - v(x_t; \theta, w^v)\right)^2\right]$$
(31)

$$= \mathbb{E}\left[\sum_{t=0}^{T} \frac{1}{2} \delta_t^2(\theta, w^v)\right].$$
(32)

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

$$\nabla_{\theta, w^{v}} \mathcal{L}^{TD} = \frac{\partial \mathcal{L}^{TD}}{\partial \delta} \cdot \frac{\partial \delta}{\partial v} \cdot \nabla_{\theta, w^{v}} v(\theta, w^{v}), \qquad (33)$$

$$= \mathbb{E}\left[\sum_{t=0}^{T} \delta_t \cdot (-1) \cdot \nabla_{\theta, w^v} v(x_t; \theta, w^v)\right], \qquad (34)$$

$$= \mathbb{E}\left[\sum_{t=0}^{T} -\nabla_{\theta^{v}} v(x_{t}; \theta, w^{v}) \cdot \delta_{t}\right].$$
(35)

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

$$\frac{\partial \delta}{\partial v} = \frac{\partial (r_{t+1} + \gamma v_{t+1} - v_t)}{\partial v_t} = -1, \qquad (36)$$

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

$$\mathcal{J}^{TD} = \mathbb{E}\left[\sum_{t=0}^{T} r_{t+1} + \gamma v_{t+1} - v_t\right] = \mathbb{E}\left[\sum_{t=0}^{T} \delta_t\right].$$
(37)

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

$$\nabla_{\theta, W^{\pi}} \mathcal{J}^{TD} = \mathbb{E} \left[\sum_{t=0}^{T} \nabla_{\theta, W^{\pi}} \log \pi(g_t | x_t; \theta, W^{\pi}) \cdot \delta_t \right].$$
(38)

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

A.4 COMBINED REWARD MAXIMIZATION OBJECTIVE FOR PLACE FIELD PARAMETERS

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. 27, and the critic employs gradient descent for Eq. 35. 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

$$\nabla_{W^{\pi},w^{\upsilon},\theta}\mathcal{J} = \nabla_{W^{\pi},w^{\upsilon},\theta}\mathcal{J}^{TD} - \nabla_{W^{\pi},w^{\upsilon},\theta}\mathcal{L}^{TD}$$
(39)

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

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

$$= \mathbb{E}\left[\sum_{t=0}^{T} \left(\nabla_{W^{\pi}, w^{v}, \theta} \log \pi(g_{t} | x_{t}; W^{\pi}, \theta) + \nabla_{W^{\pi}, w^{v}, \theta} v(x_{t}; w^{v}, \theta)\right) \delta_{t}\right].$$
(42)

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. 27, which is equivalent to the negative log likelihood

$$\nabla_{W^{\pi},w^{v},\theta}\mathcal{L} = -\nabla_{W^{\pi},w^{v},\theta}\mathcal{J}^{TD} + \nabla_{W^{\pi},w^{v},\theta}\mathcal{L}^{TD}$$
(43)

$$= -\mathbb{E}\left[\sum_{t=0}^{T} \nabla_{W^{\pi}, w^{v}, \theta} \log \pi(g_{t}|x_{t}; W^{\pi}, \theta) \cdot \delta_{t}\right] + \mathbb{E}\left[\sum_{t=0}^{T} -\nabla_{W^{\pi}, w^{v}, \theta} \tilde{v}(x_{t}; w^{v}, \theta) \cdot \delta_{t}\right], \quad (44)$$

$$= \mathbb{E}\left[\sum_{t=0}^{T} \nabla_{W^{\pi}, w^{v}, \theta} - \log \pi(g_{t}|x_{t}; W^{\pi}, \theta) \cdot \delta_{t}\right] + \mathbb{E}\left[\sum_{t=0}^{T} - \nabla_{W^{\pi}, w^{v}, \theta} \tilde{v}(x_{t}; w^{v}, \theta) \cdot \delta_{t}\right], \quad (45)$$

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

$$\tag{46}$$

which is the same update rule used in Wang et al. (2018); Mnih et al. (2016) 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 (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. 42, 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

$$\Delta \boldsymbol{w}^{v}(t+1) = \eta \delta_{t} \boldsymbol{\phi}(x_{t}) \quad , \quad \Delta \boldsymbol{W}^{\pi}(t+1) = \eta \delta_{t} \tilde{\boldsymbol{g}}_{t} \boldsymbol{\phi}(x_{t})^{\top} , \qquad (47)$$

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

$$\Delta \boldsymbol{\theta}(t+1) = \eta_{\theta} \delta_t \left(\boldsymbol{w}_v(t) + \boldsymbol{W}_{\pi}^{\top}(t) \cdot \tilde{\boldsymbol{g}}_t \right) \nabla_{\theta} \boldsymbol{\phi}(x_t; \boldsymbol{\theta}), \qquad (48)$$

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

$$\delta_{i,t}^{bp} = \delta_t \left(w_i^v(t) + W_{ji}^\pi(t) \cdot \tilde{g}_{t,j} \right) , \qquad (49)$$

$$\Delta \alpha_{i,t} = \eta_{\alpha} \cdot \delta_{i,t}^{bp} \cdot \phi_i(x_t) \cdot \left(\frac{2}{\alpha_i}\right) , \qquad (50)$$

$$\Delta \lambda_{i,t} = \eta_{\lambda} \cdot \delta_{i,t}^{bp} \cdot \phi_i(x_t) \cdot \left(\frac{x_t - \lambda_i}{\sigma_i^2}\right), \qquad (51)$$

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

$$\Delta \sigma_{i,k,t} = \eta_{\sigma} \cdot \delta_{i,t}^{bp} \cdot \phi_{i,k}(x_t) \cdot \left(\frac{(x_t - \lambda_i)^2}{\sigma_{i,k}^3}\right), \qquad (52)$$

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.

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 $\eta_{\theta} = 0.0001$.

The critic estimates the value as

$$v(x_t) = \sum_{i=1}^{N} w_i \phi_i(x_t, \boldsymbol{\theta}_i), \qquad (53)$$

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

$$\phi_i(x_t, \boldsymbol{\theta}_i) = \alpha_i^2 \exp\left(-\frac{1}{2\sigma_i^2}(x_t - \lambda_i)^2\right).$$
(54)

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

$$\frac{d}{dt}\boldsymbol{w}(t) = \boldsymbol{M}(t)(\boldsymbol{w}^{V} - \boldsymbol{w}(t)), \qquad (55)$$

$$\frac{d}{dt}\boldsymbol{\theta}_i(t) = \epsilon \, w_i(t) \mathbb{E}_{x_t} \nabla_{\boldsymbol{\theta}_i} \phi_i(x_t, \boldsymbol{\theta}_i) \delta_t \,, \tag{56}$$

The key assumption we make is that the dimensionless ratio of learning rates, ϵ is perturbatively small

$$\epsilon = \frac{\eta_{\theta}}{\eta} \ll 1,\tag{57}$$

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

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

$$\boldsymbol{\theta}(t) = \boldsymbol{\theta}_0(t) + \epsilon \boldsymbol{\theta}_1(t) + \epsilon^2 \boldsymbol{\theta}_2(t) + \dots$$
 (58)

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

The leading zeroth order dynamics are

$$\frac{d}{dt}\boldsymbol{\theta}_0(t) = 0, \ \frac{d}{dt}\boldsymbol{w}_0(t) = \boldsymbol{M}_0(\boldsymbol{w}_V - \boldsymbol{w}_0(t)),$$
(59)

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

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

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

$$\frac{d}{dt}\alpha_i(t) = \epsilon \, \frac{2}{\alpha_i(t)} w_i \sum_j M_{ij}(w_j^V - w_j) \,, \tag{61}$$

The zero-th order solution to Eq. 55 is

$$\Delta \boldsymbol{w}_0(t) \equiv \boldsymbol{w}_V - \boldsymbol{w}_0(t) = \exp\left(-\boldsymbol{M}t\right) \boldsymbol{w}_V, \tag{62}$$

$$\boldsymbol{w}_0(t) = [\boldsymbol{I} - \exp\left(-\boldsymbol{M}t\right)]\boldsymbol{w}_V, \qquad (63)$$

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

$$\frac{d}{dt}\boldsymbol{\alpha}_{1}(t) = 2\boldsymbol{\alpha}_{0}^{-1} \odot [\boldsymbol{I} - \exp\left(-\boldsymbol{M}t\right)] \boldsymbol{w}_{V} \odot \boldsymbol{M} \exp\left(-\boldsymbol{M}t\right) \boldsymbol{w}_{V}.$$
(64)

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

$$\boldsymbol{\alpha}_{1}(t) = 2\sum_{k\ell} (\boldsymbol{w}_{V} \cdot \boldsymbol{u}_{k}) (\boldsymbol{u}_{\ell} \cdot \boldsymbol{w}_{V}) (\boldsymbol{u}_{k} \odot \boldsymbol{u}_{\ell}) \left[(1 - e^{-\lambda_{k}t}) - \frac{\lambda_{k}}{\lambda_{k} + \lambda_{\ell}} (1 - e^{-(\lambda_{k} + \lambda_{\ell})t}) \right], \quad (65)$$

we can approximate this at late times as

$$\lim_{t \to \infty} \boldsymbol{\alpha}_1(t) \approx 2\boldsymbol{w}_V \odot \boldsymbol{w}_V \,. \tag{66}$$

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 are heavily involved in the reproduction of the value function are upweighted in their amplitude.

B.2 FIELD CENTER

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

$$\frac{d}{dt}\lambda_i(t) \approx \epsilon \, \frac{x_t - \lambda_i}{\sigma_i^2} \, w_i \phi_i(x) \sum_j \phi_j(x) (w_j^V - w_j) \,. \tag{67}$$

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

$$\left\langle \frac{(x_t - \lambda_i)}{\sigma^2} \phi_i(x) \phi_j(x) \right\rangle \approx \frac{\mu_{ij} - \lambda_i}{\sigma^2} M_{ij} ,$$
 (68)

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

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

Following the solution in Eq. 63, we can approximate this at late times as

$$\lim_{t \to \infty} \boldsymbol{\lambda}(t) - \boldsymbol{\lambda}(0) \approx \epsilon \boldsymbol{w}^{V} \odot \left[\left(\frac{2}{\boldsymbol{\sigma}^{2}} + \frac{1}{\sigma_{x}^{2}} \right)^{-1} \left(\frac{1}{\boldsymbol{\sigma}^{2}} (\boldsymbol{\lambda}(0) + \bar{\boldsymbol{\lambda}}) + \frac{1}{\sigma_{x}^{2}} \bar{\mu}_{x} \right) - \boldsymbol{\lambda}(0) \right] \odot \boldsymbol{w}^{V}.$$
(70)

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

$$\lambda_i(t) - \lambda_i(0) \approx \frac{\eta_\lambda}{\eta} \left(\frac{2}{\sigma_i^2} + \frac{1}{\sigma_x^2}\right)^{-1} \left[\frac{\bar{\lambda} - \lambda_i(0)}{\sigma_i^2} + \frac{\bar{\mu}_x - \lambda_i(0)}{\sigma_x^2}\right] w_{v,i}^2(t) , \ \eta_\lambda \ll \eta , \tag{71}$$

where λ 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.



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 SUCCESSOR REPRESENTATION AGENT

The generalized temporal difference error is given by

$$\delta_{t,j}^{SR} = \phi_j(x_t) + \gamma \psi_j^{\pi}(x_{t+1}) - \psi_j^{\pi}(x_t), \qquad (72)$$

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

$$\psi_i^{\pi}(x_t) = \sum_{i}^{N} [U_{ji}]_+ \phi_i(x_t) , \qquad (73)$$

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.

$$\Delta U_t = \phi_i(x_t) \cdot \delta_{t,j}^{\top}, \qquad (74)$$

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).

D METRIC REPRESENTATION LEARNING OBJECTIVE

The hippocampus is known to learn and represent spatial maps even in the absence of rewards, enabling rapid navigation to new locations when required (Tolman, 1948; Steele & Morris, 1999; Tse et al., 2007). This requires reorganization of place fields in non-rewarded conditions, which has been proposed as a mechanism for learning a predictive map that estimates future spatial occupancy (Mehta et al., 1997; Stachenfeld et al., 2017). To describe this non-reward-based reorganization, the successor representation algorithm (Dayan, 1993) has been used. More recently, an auxiliary predictive objective has been proposed (Fang & Stachenfeld, 2023).

Here, we present a simple predictive objective for place field reorganization that is independent of rewards. We introduce a previously described objective called the Metric Representation (MR), which learns a low-dimensional representation of an environment using place field activity and a biologically plausible learning rule that is modulated by a path integration-derived Temporal Difference error. This representation allows an agent to predict its current coordinates $z(x_t)$ and perform vector subtraction to rapidly navigate to recalled goals (Foster et al., 2000; Kumar et al., 2024b). However, representation learning was not studied using this objective. Recently, a similar objective was proposed to learn a spatial map using local learning rules, although as a high-dimensional representation (Stöckl et al., 2024).

The dimensionality of the coordinate prediction $z(x_t)$ is equal to the dimensionality of the environment, calculated through a linear summation of place field activity:

$$z_j(x_t) = \sum_i^N W_{ji}^{MR} \phi_i(x_t), \quad \boldsymbol{z} \in \mathbb{R}^D.$$
(75)

When the agent accurately predicts its coordinates in the environment, the following path integration-derived self-consistency equation holds:

$$z_j(x_{t+1}) \equiv z_j(x_t) + a_j(x_t),$$
(76)

$$z_j(x_{t+1}) - z_j(x_t) - a_j(x_t) \equiv 0, \qquad (77)$$

where $a_j(x_t)$ is the true displacement of the agent in the environment. However, if the prediction is inaccurate, Eq. 77 can be reformulated into a temporal difference error for each dimension j of the environment as described by Foster et al. (2000); Kumar et al. (2024b):

$$\boldsymbol{\chi}_t = z_j(x_{t+1}) - z_j(x_t) - a_j(x_t), \quad \boldsymbol{\chi} \in \mathbb{R}^D.$$
(78)

This one step prediction error (χ_t) can be expressed as a loss function, similar to Fang & Stachenfeld (2023) without the temporal discounting factor:

$$\mathcal{L}^{MR} = \mathbb{E}_{g \sim \pi} \left[\sum_{t=0}^{T} \frac{1}{2} \boldsymbol{\chi}_{t}^{2} \right] = \mathbb{E} \left[\sum_{t=0}^{T} \frac{1}{2} (\boldsymbol{z}(x_{t+1}; W^{MR}) - \boldsymbol{z}(x_{t}; W^{MR}) - \boldsymbol{a}(x_{t}))^{2} \right], \quad (79)$$

which can be minimized by gradient descent by optimizing both the coordinate readout weights (W_{MR}) and place field parameters $(\theta \in \{\alpha, \lambda, \sigma\})$:

$$\nabla_{W^{MR},\theta} = \mathbb{E}\left[\sum_{t=0}^{T} \boldsymbol{\chi}_t \nabla \boldsymbol{\phi}(\boldsymbol{x}_t;\theta)^\top\right].$$
(80)

The gradient updates were implemented in an online manner:

$$\Delta \mathbf{W}^{MR}(t+1) = \eta \boldsymbol{\chi}_t \boldsymbol{\phi}(x_t)^{\top}, \qquad (81)$$

$$\Delta\theta(t+1) = \eta_{\theta} \mathbf{W}_{MR}^{\dagger}(t) \boldsymbol{\chi}_{t} \nabla_{\theta} \boldsymbol{\phi}(x_{t};\theta), \qquad (82)$$

We can analyze how the different temporal difference residues (both the canonical reward-dependent and newly proposed metric representation-based) influence place field reorganization and agent policy learning performance by propagating the residues through a combination of actor, critic, and metric representation weights:

$$\boldsymbol{\delta}_{t}^{bp} = \delta_{t} \left(\beta_{v} \boldsymbol{w}_{v}(t) + \beta_{\pi} \boldsymbol{W}_{\pi}^{\top}(t) \cdot \tilde{\boldsymbol{g}}_{t} \right) + \beta_{MR} \mathbf{W}_{MR}^{\top}(t) \boldsymbol{\chi}_{t} , \qquad (83)$$

$$\Delta \theta_i(t+1) = \eta_\theta \boldsymbol{\delta}_t^{bp} \nabla_\theta \boldsymbol{\phi}(x_t; \theta) \,. \tag{84}$$

This can be done by setting the weighting of each component $\beta_v, \beta_\pi, \beta_{MR} \in \{0, 1\}$. Refer to Fig. 15 for policy convergence performance in both the 1D and 2D environments when using different combinations to learn place field representations.

E 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 ($\theta \in \{\alpha, \lambda, \sigma, w^v, W^\pi\}$).

$$\theta_{t+1} = \theta_t + \xi_t \,, \tag{85}$$

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 Fig. 7.

F SUPPLEMENTARY FIGURES



Supplementary Figure 1: Influence of place field parameter optimization for a single seed. Example change in individual field's spatial selectivity ($\phi(x)$, colored), mean firing activity at a location $(\sum_{i=1}^{N} \phi_i(x))$, field density which is the number of Center of Mass (COM) in a location after smoothing with a Gaussian kernel density estimate (gKDE) (gKDE(COM), blue) and, the frequency of being in a location $(p_{RM}(x))$, when optimizing different combinations of field parameters $(\alpha, \lambda, \sigma)$ during reward maximization (RM). The location in which the highest value for mean firing activity, field density and frequency is attained is indicated by a red, blue and black vertical dash line respectively. Optimizing a (A) small number (N = 16) and (B) large number of place fields yields a similar high mean firing rate at the reward location followed by the start location. However, the field density evolves differently when in the low field regime, (A) a high density emerges at the reward location in the early stages of learning, but it shifts to the start location at later stages of learning. This effect was observed in a recent experiment where place fields which initially encoded the reward location, gradually shifted backward towards the corresponding start location. This shift led to a decrease in place fields specifically coding for the reward, suggesting that the hippocampal representation reorganizes to predictively code for the reward (Yaghoubi et al., 2024). Whether experiments demonstrate such misalignment between place field density and mean firing rate needs to be analyzed. Based on the ablation studies (Fig. 4A,B), mean firing rate will be a stronger indicator of learning performance than field density. (B) In the high field regime, a high field density at the reward location remains stable throughout learning. Note that COM changes only when the place field centers are optimized ($\Delta\lambda$). Distribution is shown for a single seed run for a homogeneous place field population that has been initialized by with equal spacing between field centers ($\lambda \in [-1, 1]$), equal amplitude ($\alpha = 0.5$) and width ($\sigma = 0.01$). Refer to Fig. 2 for general place field reorganization over different seeds.



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 reward location. Agents with N = 256 place fields and Gaussian noise injected to field parameters $(\sigma_{noise} = 0.0001)$ were trained to navigate to a reward location at $x_r = 0.75$ for 50,000 trials, thereafter the reward location was shifted to $x_r = -0.2$ for the next 50,000 trials. (A) Place field density at the start of learning was uniformly distributed (left) and increased near the first reward location at the end of the first 50,000 trials (center). After the shift in reward location, a high density of fields emerged at the new reward location (right). The black line shows the learned policy, where a velocity of 0.1(-0.1) indicates moving right (left). Agents learn to navigate to the reward location, both before and after the shift. (B) Example distribution of individual place fields before learning (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 to model heterogeneous place field population. Notice the backward shift of some place fields that were at the initial reward location to the new reward location. (C) About 2.6% of the place fields coding for the initial reward at $x_r = 0.75$ (green dots) shifted to the new reward location at $x_r = -0.2$ (about 19 of the 734 green dots are within the blue circle). Other place fields at $x_r = -0.2$ increased their firing rate to encode the new reward location. We see a large number of fields shifting backward, though not entirely to the new reward location. Shaded area shows 95% 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.



Supplementary Figure 4: Weak feature learning with large number of place fields. Critic w_i^v and actor W_{ii}^{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); Kumar et al. (2022); Frémaux et al. (2013). (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, $[10^{-5}, 0.1]$ 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 and MR agent architecture, and representation dynamics. (A) Successor Representation (SR) agent architecture to learn a navigational policy and the SR place fields. Only the synapses from the initialized place field (ϕ_{fixed}) to the actor (red) and critic (green), and the synapses (U) to the SR fields (ψ) were plastic. Refer to App. C for implementation details. (B) Left: Metric Representation (MR) agent architecture learns to predict the agent's coordinates in an environment. The coordinate readout is a linear summation of place field activity, and its dimension is the same as the displacement in the environment. The agent learns to predict its coordinates by minimizing a path integration derived temporal difference error χ_t . The gradient updates are performed on both the coordinate readout weights W_{ii}^{MR} and place field parameters α, λ, σ . The agent learns to navigate to the reward location only by updating the actor and critic weights, without influencing place field parameters Refer to App. D for details. Hence, place fields in the MR agent will reorganize even in the absence of rewards. Right: Change in MR agent's coordinate estimation in a 1D track across trials (T = 0, 1, 10, 50000). Coordinate estimation was close to zero during W_{ii}^{MR} initialization. After 10 trial, the agent starts to show a monotonic increase in coordinate estimation as the agent moves from x = -1 to x = 1. By 50,000 trials, the agent's coordinate estimation becomes stable. (C) Average change for 16 and 64 place fields' size (firing rate greater than 10^{-3} in the track) (top row) and center of mass (bottom row) when SR, RM and MR agents navigate in a 1D track with the absolute change reflected in the y axis. Shaded area shows 95% CI over 5 seeds for agents with 16 and 64 place fields. (D) Spatial representation similarity matrix for SR (top row) and RM (middle row) and MR (bottom row) agents in a 1D track is visualized by taking the dot product of the place field activity at each location. (E) Difference in correlation between the proportion of time spent in a location between SR, RM, MR agents. (F) The correlation between the individual field firing rates learned by SR, RM, and MR agents rapidly diverge but remain positively correlated. (G) The correlation between the spatial representation similarity matrices (purple) learned by SR, RM and MR rapidly diverge in the early learning phase but stabilize and remain positively correlated in later phases.



Supplementary Figure 6: Field elongation in 2D arena. (A-C) 2D Place field distortion dynamics by SR (A), RM (B), and MR (C) agents as learning proceeds. Numbers in yellow on the obstacle indicates (Field ID)-(Maximum firing rate). (D) Average change for 256 place fields' size (top row) and center of mass () (bottom row) when SR, RM and MR agents navigate in a 2D arena with the absolute change reflected in the y axis. Area was determined by computing the firing rate that was greater than 10^{-3} in the arena. The 2D arena was divided into three sub-areas to track COM movement 1) away from the reward location, 2) the corridor from right to left, and 3) towards the start location. All three agents showed an increase in field area and backward COM shift towards the start location. Shaded area shows 95% CI over 3 seeds. (E) Change in coordinate readout weights in a 2D environment. Each plot indicates the synaptic weights W_{ji}^{MR} from place fields to the x (top row) and y dimensions of the 2D environment respectively. Weights were randomly initialized in trial 0. As the agent explores the environment, the weights converge to reflect a spatial map where the coordinate estimation for the X and Y axes increase monotonically when the agent moves left to right and bottom to top respectively, similar to Foster et al. (2000); Kumar et al. (2024b) which used a similar path-integration TD error but with eligibility traces instead.



Supplementary Figure 7: Noise amplitude monotonically influences population vector correlation 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_v) 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 performance. 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) , 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 representation. (A) Reward maximization performance variability increases when noise magnitude increases. (B) With no noise injection, variance in parameter update is initially positively correlated with field amplitude (blue). When a small amount of noise is added, fields with a larger mean amplitude show a smaller variance in change in parameter while fields with a smaller amplitude show higher variance. Conversely, when the magnitude of noise is further increased (purple), fields with a higher 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 critic weights) is high and positively correlated when the noise magnitude is low. This correlation decreases and becomes weakly positive when $\sigma_{noise} = 0.001$. This supports the claim that in the low noise regime, fields with a high amplitude are more involved in policy learning and hence drift less or are more stable to maintain performance integrity. (D) Population vector correlation decreases at a faster rate than the similarity matrix when noise magnitude increases. (E) Representation similarity correlation decreases as the noise magnitude increases, but at a slower rate than PV correlation. (F) Proportion of fields that are active (average fraction of fields with firing rate less than 0.05, 0.1,0,25) continues to increase with higher noise magnitude. (G) Introducing Gaussian noise 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. E 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 or lose selectivity, and most fields increased in size (σ^2) while some did not. The first two were observed by Qin et al. (2023) who analyzed Gonzalez et al. (2019). Each color corresponds to the dynamics of a specific field, with 5 example fields shown.



Supplementary Figure 9: Noisy place field parameter update replicates drift dynamics seen in neural data. (A) Place field centroids becomes distinctively different across trials, after stable navigation performance was attained at trial 25,000, similar to Ziv et al. (2013); de Snoo et al. (2023). Each place field's centroid position was sorted according to trial 25,000, 125,000 and 195,000. (B) When no Gaussian noise is added to place field parameters $(\alpha, \lambda, \sigma)$, place field optimization alone does not cause centroids to shift. Instead, adding small Gaussian noise $(\sigma_{noise} \in \{0.0001, 0.00025, 0.0005\})$ replicates the gradual shift in centroid position across trials (25,100 to 125,000) as seen in Qin et al. (2023); Ziv et al. (2013); Geva et al. (2023). When the noise magnitude is high e.g. $\sigma_{noise} \geq 0.001$, centroids shift rapidly to a new location, similar to the random shuffle or null hypothesis seen in Ziv et al. (2013); Qin et al. (2023); Geva et al. (2023). (A-B) Analysis was done for 64 place fields aggregated over 10 agents initialized with different seeds to have 640 fields in total. (C) Example graph topology for one agent with N = 64 place fields with Gaussian noise $\sigma_{noise} = 0.00025$ added to field parameters. Each node indicates a place field's centroid position across learning, and the edge is weighted by the normalized (between 0 to 1) cosine distance between each node that is less than 0.55. Red, green, blue, orange, black nodes indicate centroids initialized at the reward, start, end of track near the reward, end of track near the start locations and the middle of the track respectively. As learning progressed, the cosine distance between each centroid changed 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 phenomenon where a high density of fields emerge at the reward and start locations.



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 an obstacle. (A) Agents started at the same location $x_{start} = (0.0, 0.75)$ and had to navigate to a target that changed to a new location every 50,000 trials following the sequence $(x_r \in [(0.75, -0.75), (-0.75, 0.75), (0.75, 0.75), (-0.75, -0.75)])$. Increasing the noise magnitude 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 $(x_{min} = -0.2, x_{max} = 0.2, y_{min} = -1.0, y_{max} = 0.5)$ for the first 50,000 trials. After which, the location of the obstacle was shifted up to $(x_{min} = -0.2, x_{max} = 0.2, y_{min} = -0.5, y_{max} = 1.0)$ while the start and target location was the same. Agents with a noise magnitude $\sigma_{noise} = 0.00025$ showed the highest average reward maximization performance followed by $\sigma_{noise} = 0.0005$. A high noise magnitude ($\sigma_{noise} = 0.001$) disrupted learning performance while agents without noisy field updates ($\sigma_{noise} = 0.0$) did not learn to navigate around the new obstacle. Note that field amplitudes and widths were clipped to be between $[10^{-5}, 2]$ and $[10^{-5}, 0.5]$ respectively to ensure the Σ covariance matrix in 2D place fields remained valid for matrix inversion. Performance was averaged over agents initialized with different number of 2D place fields ($N \in \{64, 144, 256, 576\}$) with the diagonals of the field width initialized with $\Sigma = 0.01$ and constant amplitude $\alpha = 1.0$, over 30 different seeds. Shaded area is 95% CI.



Supplementary Figure 13: Using the same learning rates for the place field parameters and actor-critic recovers the same phenomena of a high field density emerging at reward location followed by the start location, and field elongation against the agent's trajectory. (A) Each place field's amplitude, center and width were sampled from a uniform distribution of $[0, 1], [-1, 1], [10^{-5}, 0.1]$ respectively to model heterogeneous place field distribution. After learning, a high density (number) of fields emerged at the start (green dash) and reward (red area) location, similar to Fig. 1B (right) and Sup. Fig. 2B. This phenomenon is consistent across different numbers of place fields. Shaded area is 95% CI over 50 different seeds. (C) In a 2D arena with obstacles, place fields elongate from the reward location (red circle) back to the start location (green circle), while narrowing along the corridor with an obstacle (gray), similar to Fig. 2F. Learning rates for the actor, critic and place field parameters were $\eta = \eta_{\theta} = 0.0005$.



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 to optimize place field parameters. We evaluated the speed of policy learning when optimizing (A) heterogeneously distributed place field population in the 1D track and (B) homogeneously distributed place field population in the 2D arena using: (1) fixed place field parameters (blue), (2) backpropagating the TD error δ_t through the actor weights $W_{\pi}^{\dagger} \tilde{g}_t$ $(\beta_{\pi} = 1, \beta_{v} = 0, \beta_{MR} = 0, \text{ orange}), (3)$ backpropagating the TD error through the critic weights w^{v} ($\beta_{\pi} = 0, \beta_{v} = 1, \beta_{MR} = 0$, green), (4) backpropagating the path integration derived TD error χ_t through the metric representation weights W_{MR} ($\beta_{\pi} = 0, \beta_{\nu} = 0, \beta_{MR} = 1$, red) while learning the value function and policy by optimizing only the readout critic and actor weights, (5) backpropagating the TD error through both the actor and critic weights, otherwise called the Reward Maximization agent ($\beta_{\pi} = 1, \beta_{v} = 1, \beta_{MR} = 0$, purple), (6) backpropagating the TD error through both the actor and critic weights and the path integration based TD error through the metric representation weights, ($\beta_{\pi} = 1, \beta_{\nu} = 1, \beta_{MR} = 1$, brown). The combined RM+MR objective used for place field parameter optimization achieved the fastest policy learning, similar to Stachenfeld et al. (2017) when the number of fields was low ($N = \{4, 8, 16, 32\}$ in 1D and $N = \{4, 8\}$ in 2D). With more fields, the reward maximization agent (RM, purple) was almost as effective as the combined objective (RM + MR, brown). Optimizing place field parameters using only the actor weights led to the slowest policy convergence (orange), nevertheless faster than using fixed place fields. The same learning rates were used when the number of fields were increased. Hence, tuning learning rates should improve the stability of policy learning, especially in the 2D arena for the agent with the combined RM+MR objective. Shaded area indicates 95% CI over 50 random seeds.