

# Emergent Elasticity in the Neural Code for Space

Samuel Ocko<sup>a</sup>, Kiah Hardcastle<sup>b</sup>, Lisa Giocomo<sup>b</sup>, and Surya Ganguli<sup>a,b</sup>

<sup>a</sup>Department of Applied Physics, Stanford University, Stanford, CA 94305

<sup>b</sup>Department of Neurobiology, Stanford University, Stanford, CA 94305

**Upon encountering a novel environment, an animal must construct a consistent environmental map, as well as an internal estimate of its position within that map, by combining information from two distinct sources: self-motion cues and sensory landmark cues. How do known aspects of neural circuit dynamics and synaptic plasticity conspire to accomplish this feat? Here we show analytically how a neural attractor model that combines path integration of self-motion cues with Hebbian plasticity in synaptic weights from landmark cells can self-organize a consistent map of space as the animal explores an environment. Intriguingly, the emergence of this map can be understood as an elastic relaxation process between landmark cells mediated by the attractor network. Moreover, our model makes several experimentally testable predictions, including: (1) systematic path-dependent shifts in the firing field of grid cells towards the most recently encountered landmark, even in a fully learned environment, (2) systematic deformations in the firing fields of grid cells in irregular environments, akin to elastic deformations of solids forced into irregular containers, and (3) the creation of topological defects in grid cell firing patterns through specific environmental manipulations. Taken together, our results conceptually link known aspects of neurons and synapses to an emergent solution of a fundamental computational problem in navigation, while providing a unified account of disparate experimental observations.**

Grid Cells | Border Cells | Theoretical Neuroscience

Correspondence: [socko@stanford.edu](mailto:socko@stanford.edu)

How might neural circuits learn to create a long term map of a novel environment and use this map to infer where one is within the environment? This pair of problems are challenging because of their nested, chicken and egg nature. To localize where one is in an environment, one first needs a map of the environment. However, in a novel environment, no such map is yet available, so localization is not possible. Instead, neural circuits must create a map over time, through exploration in a novel environment, without initially having access to any global estimate of position within the environment. This chicken and egg problem is known in the robotics literature as Simultaneous Localization and Mapping (SLAM) (1). Here we explore how known aspects of neural circuit

dynamics and synaptic plasticity can conspire to self-organize, through exploration, a neural circuit solution to the problem of creating a global, consistent map of a novel environment. In particular, neural circuits receive two fundamentally distinct sources of information about position: (1) signals indicating the speed and direction of the animal, which can be path-integrated over time to update the animal's internal estimate of position, and (2) sensory cues from salient, fixed landmarks in the environment. To create a map of the environment, neural circuits must combine these two distinct information sources in a self-consistent fashion so that sensory cues and self-motion cues are always in co-register.

For example, consider the act of walking from landmark A to landmark B. Sensory perception of landmark A triggers a pattern of neural activity, and subsequent walking from A to B evolves this activity pattern, through path integration, to a final pattern. Conversely, sensory perception of landmark B itself triggers a neural activity pattern. Any circuit that maps space must obey a fundamental self-consistency condition: the neural activity pattern generated by perception of A, followed by path integration from A to B, must match the neural activity pattern triggered by perception of B alone. Only in this manner can neural activity patterns be in one to one correspondence with physical positions in space, and become independent of the past trajectory used to reach any physical location.

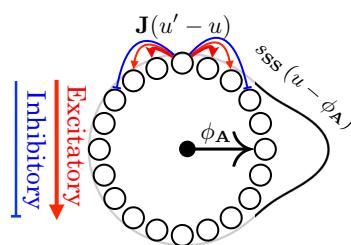
In the following, we develop an analytic theory for how neuronal dynamics and synaptic plasticity can conspire to self-organize such a self-consistent neural map of space upon exploration through a novel environment. Moreover, our analytic theory makes experimentally testable predictions about neural correlates of space. Indeed, many decades of recordings in multiple brain regions have revealed diverse neural correlates of spatial maps in the brain. In particular, the medial entorhinal cortex (MEC) contains neurons encoding for direction, velocity, landmarks, as well as grid cells exhibiting striking firing patterns reflecting an animal's spatial location (2–6). Moreover, the geometry of these firing patterns depends

on the shape of the environment being explored (7–10). In particular, these grid firing patterns can be deformed in irregular environments (11, 12), in a manner evocative of deformations of solids forced into an irregular container, suggesting a mechanical model for these deformations (13–15). Also, these firing patterns are not simply driven by current sensory cues; there is evidence for path integration (16–18) in that firing patterns appear almost immediately (2), phase differences are preserved across environments (19), firing patterns become noisier the longer an animal has spent away from a landmark (20, 21), and can be shifted depending on which landmark the animal has most recently encountered (22, 23).

Despite this wealth of experimental observations, no mechanistic circuit model currently explains how known aspects of neuronal dynamics and synaptic plasticity can conspire to learn, through exploration, a self-consistent internal map of a novel environment that both behaves like a deformable medium, and also retains, at higher-order, some knowledge of recently encountered landmarks. Here, we show how an attractor network that combines path integration of velocity with Hebbian learning (22, 24, 25) of synaptic weights from landmark cells, can self-organize to generate all of these outcomes. Intriguingly, a low dimensional reduced model of the combined neuronal and synaptic dynamics provides analytical insight into how self-consistent maps of the environment can arise through an emergent, elastic relaxation process involving the synaptic weights of landmark cells.

## Model reduction of an attractor network coupled to sensorimotor inputs

Our theoretical framework assumes the existence of three interacting neural components: (1) an attractor network capable of realizing a manifold of stable neural activity patterns, (2) a population of velocity-tuned cells that carry information about the animal's motion, and (3) a population of sensory driven landmark cells that fire if and only if the animal is in a particular region of space. Our goal will be to understand how the three populations can interact together and self-organize through synaptic plasticity, sculpted by experience, to create a self-consistent internal map of the environment. Here, we describe the neuronal and synaptic dynamics of each component in turn, as well as describe a model reduction approach to obtain a low dimensional reduced description of the entire plastic circuit dynamics. Our low dimensional description provides insight into how self-consistency of the neural map emerges naturally through an elastic relaxation process between landmarks.



**Fig. 1.** Schematic of a ring attractor with short-range excitation (red arrows) and longer range inhibition (blue arrows). This yields a 1D family of bump-attractor states  $s_{SS}(u - \phi_A)$ , which are mapped onto a single periodic variable  $\phi_A$  representing the peak of the bump pattern.

## A manifold of stable states from attractor network dynamics

We first consider a one-dimensional attractor network consisting of a large population of neurons whose connectivity is determined by their position on an abstract ring, as in Fig. 1. For analytical simplicity, we take a neural field approach (26), so that position on the ring of neurons is described by a continuous coordinate  $u$ , with the firing rate of a neuron at position  $u$  given by  $s(u)$ . Each neuron interacts with neighboring neurons through a translation invariant connectivity, yielding the dynamics

$$\frac{ds(u)}{dt} = -s(u) + \mathcal{F} \left[ \int_{u'} \mathbf{J}(u - u') s(u') \right]. \quad (1)$$

Here  $\mathbf{J}(u - u')$  defines the synaptic weight from a cell at position  $u'$  to a cell at  $u$ , and  $\mathcal{F}$  is a nonlinearity. We will refer to these dynamics as  $ds/dt = \text{Dyn}[s]$ . Many appropriate choices of  $\mathbf{J}$  and  $\mathcal{F}$ , corresponding for example to short range excitation and long range inhibition, will yield a family of stable, or steady state, localized bump activity patterns  $s_{SS}(u - \phi_A)$ , parameterized by the position of their peak  $\phi_A$  (27, 28). This one-dimensional family of stable bump activity patterns can itself be thought of as ring of stable firing patterns in the space of all possible firing patterns. Just as  $u$  indexes a family of neurons on the neural sheet, the coordinate  $\phi_A$  indexes the different stable neural activity patterns, with a particular value of  $\phi_A$  corresponding to a stable bump on the neural ring centered at coordinate  $u = \phi_A$ . For simplicity we set units such that the coordinate  $u$  along the neural ring, and the coordinate  $\phi_A$  along the ring of stable attractor patterns are both periodic variables defined modulo  $2\pi$ . Thus  $u$  and  $\phi_A$  are phase variables denoting position along the neural ring and ring of bump attractor patterns respectively.

**Motions along the attractor manifold due to external inputs** So far, the attractor network described above has a ring of stable bump activity patterns parameterized by the periodic coordinate  $\phi_A$ , but these neural activity patterns are as yet unanchored to physical space. We will

eventually show how to anchor the coordinate  $\phi_A$  along the attractor manifold to the actual position of the animal in physical space. However, in order to appropriately form such an internal map of position, and thereby map the environment, the attractor state must be influenced by external inputs from both velocity and landmark sensitive cells in a self-consistent manner. We first derive a reduced description for how a general external feedforward input to the attractor network modifies its dynamics.

Suppose the attractor network is at one of its steady state bump patterns  $s_{SS}(u - \phi_A)$  centered at  $u = \phi_A$ . Further suppose that each neuron at position  $u$  on the neural ring experiences an external additive input current  $\epsilon \text{Pert}(u - \phi^P)$  that is centered, or localized on the neural ring around some other location  $u = \phi^P$ . The neural dynamics of the attractor in Eq. 1, in response to this additive external input is then modified to:

$$ds/dt = \text{Dyn}[s(u)] + \epsilon \text{Pert}(u - \phi^P).$$

When  $\epsilon$  is small, the external inputs are weak relative to the recurrent inputs that determine the shape of the bump pattern. In this situation, the evolving firing rates will be confined to the 1D manifold of steady states  $s_{SS}(u - \phi_A)$ . In essence, a small excitatory perturbation  $\text{Pert}(u - \phi^P)$  centered at position  $u = \phi^P$  on the neural ring, will translate the stable bump pattern towards the perturbation, *without* changing its shape. Therefore, to track the entire dynamics of the network, we do not need to track the firing rate of every neuron; we need only track the time dependent position of the peak of the activity bump,  $\phi_A(t)$ . Thus we can reduce the entire high dimensional neural dynamics to a low dimensional effective scalar dynamics:

$$d\phi_A/dt = \epsilon \text{Force}_A(\phi^P - \phi_A). \quad (2)$$

In App. A, we show how to analytically compute the force law  $\text{Force}_A(\phi^P - \phi_A)$  governing the velocity of the bump peak, in terms of the shape of the bump  $s_{SS}(u - \phi_A)$  and the shape of the additive input perturbation  $\text{Pert}(u - \phi^P)$ . However, for the particular external inputs we consider below, the qualitative structure of the force law as a function of the input perturbation will be highly intuitive.

**Path integration through conjunctive position velocity inputs** Following (27, 28), we achieve path integration by coupling the attractor network to conjunctive position and velocity-tuned cells such that east (west) movement-selective cells form feedforward synapses into the attractor network that are shifted in the positive (negative)  $u$  direction (Fig. 2A, B). We can use our model reduction framework via Eq. 2 to show analytically (App. B) that

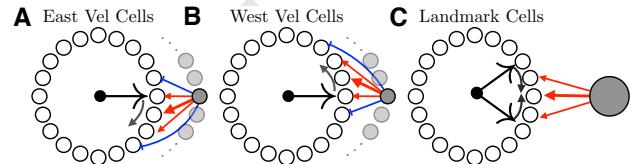
this choice of connectivity leads to path integration:

$$d\phi_A/dt = v_x \cdot k_x. \quad (3)$$

Here,  $k_x$  is a constant of proportionality that relates animal velocity to the rate of phase advance in the attractor network ( $k_x = 2\pi/\text{Field Spacing}$ ). Solving Eq. 3 allows us to recover path integration where the resulting integrated attractor phase is *only* a function of current position  $\mathbf{r}(t)$ :

$$\phi_A(t) = \phi_A(\mathbf{r}(t)) = \phi_A(0) + k_x \cdot [\mathbf{r}_x(t) - \mathbf{r}_x(0)].$$

Thus the connectivity of the conjunctive position velocity cells in Fig. 2A, B ensure that as the mouse moves east (west) along a 1D track, the attractor phase moves clockwise (counterclockwise), at a speed proportional to velocity. The collection of neurons in the attractor then trace out periodic firing patterns as a function of spatial position, all with the same period but different phases.



**Fig. 2.** **A)** When the animal moves east, east-conjunctive cells with biased outgoing connections move the attractor pattern in the positive  $u$  direction. **B)** When the animal moves west, the attractor pattern is moved in the negative  $u$  direction. **C)** Schematic of a landmark cell correcting the the attractor bump (Eq. 4). A single landmark cell will pull the peak of the bump pattern towards the peak of its efferent synaptic strength profile.

However, even though these 1D grid cell firing patterns are now a function of physical space, they still are not yet anchored to the environment. There is as yet no mechanism to set the phase of each cell relative to landmarks, and indeed these grid patterns rapidly decohere without anchoring to landmarks, as demonstrated experimentally (20, 29). Coupling the attractor network to landmark-sensitive cells can solve this problem.

**Landmark Cells** We model each landmark cell  $i$  as purely sensory driven cell with a firing rate that depends on location through  $\text{Firing}_i(t) = H_i(\mathbf{r}(t))$ . Here  $H_i(\mathbf{r})$  is the firing field of the landmark cell. An example of a landmark cell could, for example, be an entorhinal border cell (4). Every landmark cell forms feed-forward connections onto each cell in the attractor network at ring position  $u$  with a synaptic strength  $W_i(u)$ .

Consider for example a single landmark cell whose synaptic strength  $W(u)$  as a function of position  $u$  on the neural ring consists of a single bump centered at a particular location  $u = \phi^L$  (Fig. 2C). Through our model reduction framework of Eq. 2, if this landmark cell fires, then it

will exert a force on any other attractor bump pattern  $s_{\text{SS}}(u - \phi_A)$  centered at  $u = \phi_A$ , through:

$$\frac{d\phi_A}{dt} = \omega \text{Force}_A(\phi^L - \phi_A). \quad (4)$$

Here we have introduced  $\omega$  as a parameter that controls how strongly landmark cells influence the attractor phase. In essence, when each landmark cell fires, it forces the attractor state  $\phi_A$  to flow towards the phase  $\phi^L$  corresponding to the location of its maximal outgoing synaptic strength. An attractor phase  $\phi_A$  that is smaller (larger) than the landmark cell synapses' peak location  $\phi^L$  will increase (decrease) and settle down at  $\phi^L$  (Fig. 2C). Indeed for general synaptic strength patterns peaked at  $u = \phi^L$ , the force law will have the same qualitative features as  $\text{Force}_A(\phi^L - \phi_A) = \sin(\phi^L - \phi_A)$  (App. B.1).

However there is, as of yet, no mechanism to enforce consistency between the path integration dynamics of the attractor network in the absence of landmarks in Eq. 3, with the driving force exerted by a landmark cell to a particular phase  $\phi^L$  through Eq. 4. We next introduce Hebbian plasticity of efferent landmark cell synapses during exploration while *both* path integration and landmark cells are active. We then show how such plasticity yields a precise mechanism for the self-consistency required of any spatial map forming circuit.

**Hebbian Learning Between Landmark Cells and Attractor Networks.** We assume that each synapse  $W_i(u)$  from a landmark cell  $i$  to an attractor cell at position  $u$  undergoes Hebbian plasticity with some weight decay, thereby learning to reinforce attractor patterns that are active when the landmark cell fires. Moreover, we assume the dynamics of plasticity varies slowly, over a timescale  $T$  that is much longer than the timescale  $t$  over which explorations occur. Then Hebbian learning drives synaptic strengths towards the long-time average of attractor states that occur during landmark cell firing through

$$dW_i(u)/dT = \langle s(u) | i \text{ Firing} \rangle - W_i(u). \quad (5)$$

Assuming the effect of landmark cells on the attractor network is strong enough to affect the position of the bump patterns, but not strong enough to change their shape, then the long term average  $\langle s(u) | i \text{ Firing} \rangle$  of attractor patterns  $s(u)$  occurring whenever the landmark cell  $i$  fires can be written as

$$\langle s(u) | i \text{ Firing} \rangle = \int_{\phi^L} s_{\text{SS}}(u - \phi^L) \Pr(\phi_A(t) = \phi^L | i \text{ Firing}).$$

Thus all that matters for determining synaptic strength is the distribution of attractor phases that occur when the

landmark cell fires.

Now, because the learning rule is linear and the landmark cell synapses only observe attractor steady states, the Hebbian weights  $W_i(u)$  can be written as a weighted superposition of the attractor bump patterns with weighting coefficients  $\tilde{W}_i(\phi^L)$ :

$$W_i(u) = \int_{\phi^L} \tilde{W}_i(\phi^L) s_{\text{SS}}(u - \phi^L). \quad (6)$$

Furthermore, inserting Eq. (6) into Eq. (5) yields the learning dynamics of the synaptic weighting coefficients (see App. A for a proof):

$$d\tilde{W}_i(\phi^L)/dT = \Pr(\phi_A(t) = \phi^L | i \text{ Firing}) - \tilde{W}_i(\phi^L). \quad (7)$$

**Combined neural and synaptic dynamics during exploration.** By combining the effect of path integration on the attractor phase  $\phi_A$  described in Eq. (3) with the effect of multiple landmark cells with arbitrary learned weights  $\tilde{W}_i(\phi^L)$ , each acting on the phase through Eq. (4), we obtain the full dynamics of attractor phase driven by both animal velocity and landmark encounters:

$$\begin{aligned} \frac{d}{dt} \phi_A &= v_x \cdot k_x \\ + \sum_i \omega_i H_i(\mathbf{r}(t)) \int_{\phi^L} \tilde{W}_i(\phi^L) \text{Force}_A(\phi^L - \phi_A). \end{aligned} \quad (8)$$

Together, Eq. 7 and Eq. 8 reflect a complex coupled dynamics between neurons and synapses. In Eq. 7 the distribution of attractor network activity patterns, or phases, drives plasticity in synapses from landmark cells to the attractor network. In turn, these synaptic weights modify the evolution of the attractor network phase via Eq. 8.

**Coupled landmark and attractor phase dynamics in the linearized model.** The learned weights of a landmark cell are composed of a distribution of attractor network states. Linearizing  $\text{Force}_A(\phi^L - \phi_A) \approx (\phi^L - \phi_A)$ , we can simplify this representation to a single variable: the weighted average of the distribution  $\theta = \int_{\phi^L} \tilde{W}_i(\phi^L) \phi^L$ , yielding a simplified equation for describing the neuronal and synaptic outcome of navigation and learning (see App. B):

$$d\phi_A/dt = k_x \cdot v_x + \sum_i \omega_i H_i(\mathbf{r}(t)) (\theta_i - \phi_A), \quad (9)$$

$$d\theta_i/dT = \langle \phi_A(t) | \text{Cell } i \text{ Firing} \rangle - \theta_i. \quad (10)$$

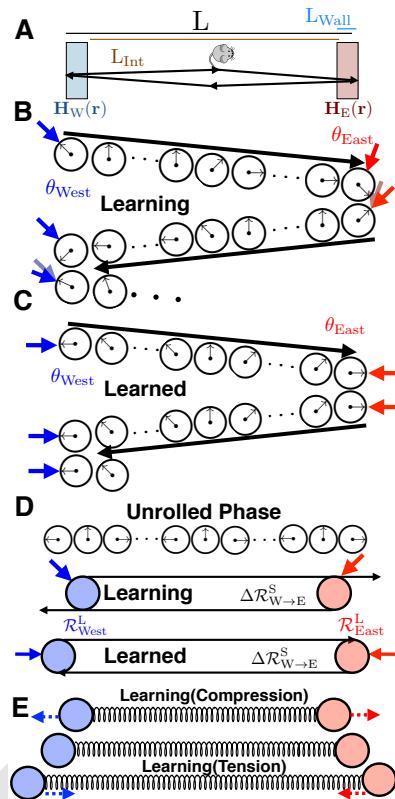
In essence Eq. 9 and Eq. 10 constitute a significant model reduction of Eq. 1 and Eq. 5. In this reduction, the entire pattern of neural activity of the attractor network is

summarized by a single number  $\phi_A$ , denoting a point, or phase, on the ring manifold of stable attractor states. Similarly, the entire pattern of synaptic weights  $W_i(u)$  from landmark cell  $i$  into the attractor network is summarized by a single number  $\theta_i$ , denoting the learned attractor network phase associated with the landmark cell's synapses. Intuitively, the reduced Eq. 9 describes both path integration and a dynamics whereby each landmark cell  $i$  attempts to *pin* the attractor phase  $\phi_A$  to the landmark cell's learned phase  $\theta_i$ , each time the physical position  $\mathbf{r}(t)$  of the animal is within the landmark's firing field  $H_i$ . In turn, synaptic plasticity described in Eq. 10 aligns the learned pinning phase  $\theta_i$  of each landmark cell  $i$  to the average of the ensemble of attractor phases  $\phi_A$  that occur when the animal is in the firing field of the landmark. As we will see below, as an animal explores its environment, this coupled dynamics between attractor phase  $\phi_A$  and landmark pinning phases  $\theta_i$  settle into a self-consistent steady state such that the attractor phase yields an internal estimate of the animal's current position that is, to first order, largely independent of the history of the animal's previous trajectory. Moreover, each landmark cell learns a pinning phase  $\theta_i$ , consistent with the location of its firing field in physical space.

## Learning a simple environmental geometry

We now examine solutions to these equations to understand how neuronal dynamics and synaptic plasticity conspire to yield a consistent map of the environment. To build intuition, first consider the linearized dynamics of Eqs. 9, 10 for the simple case of an animal moving back and forth between the walls of a 1D box of length  $L$ , at a constant speed  $v = L/\tau$ , yielding a total time of  $2\tau$  to complete a full cycle (Fig. 3A). In this environment we assume two landmark cells corresponding to the east (west) walls, with firing fields extending a distance  $L_{\text{Wall}}$  into the environment leaving an empty space  $L_{\text{Int}} = L - 2L_{\text{Wall}}$  between. Their pinning phases  $\theta_{\text{East}}$  ( $\theta_{\text{West}}$ ) encode the peak position of their outgoing synaptic weights. How does circuit plasticity yield a consistent environmental representation through exploration?

We will build intuition in the limit where  $L_{\text{Wall}} \rightarrow 0$ ,  $\omega \rightarrow \infty$ ; in this regime, landmark cells only act at the very edge, yet *fully* anchor the attractor state when the animal touches the edge. At  $t = 0$ , the animal starts at the west wall at physical position  $\mathbf{r}(0) = -L/2$ . Through Eq. 9, the west border cell pins the initial attractor phase so that  $\phi_A(0) = \theta_{\text{West}}$ . At  $t = \tau$ , the animal travels to the east wall at physical position  $\mathbf{r}(\tau) = +L/2$ , and the attractor phase advances due to path integration to become to become



**Fig. 3.** **A)** An animal moving between two landmarks at the edges of a 1D track. **B)** A single cycle of exploration as the animal moves from the west to east wall and back. When the animal encounters the west (east) wall, the attractor phase (black arrow) is pinned to the associated landmark pinning phase (blue/red arrow for west/east wall). As the animal moves from one wall to the other, the attractor phase advances from this pinned phase due to path integration. During learning, the pinning phase from any one wall, plus the phase advance due to path integration, will not equal the pinning phase of the other wall. However, plasticity will adjust the pinning phase of each wall to reduce this discrepancy (motion of red and blue arrows). During this inconsistent pre-learned state, the attractor phase at any interior position will depend on path history. **C)** After learning, the pinning phase from any one wall, plus the phase advance due to path integration, *equals* the pinning phase of the other wall, yielding a consistent internal representation of space in which the attractor phase assigned to any interior point becomes independent of path history. **D)** We can “unroll” the attractor and landmark phases into linear position variables. Thus landmark cell synapses can be thought of as points in physical space (blue and red circles). If the phase advance due to path integration *exceeds* phase difference between the pinning phases of the landmarks, then the distance between the landmark cells in unrolled phase is *closer* than the physical distance between the firing fields of the landmarks (top). Plasticity then exerts an outward force pushing the two landmark cells further apart until their separation in unrolled phase equals the physical distance between between their firing fields (bottom). **E)** In general, the changing positions in unrolled phase associated with landmark cell synapses due to synaptic plasticity can be described by a damped spring-like interaction as in Eq. 11 and Eq. 12. If the separation between the two landmark cell synapses in unrolled phase is smaller (larger) than the physical separation between their firing fields, then the spring will be compressed (extended), yielding an outward (inward) force. This force will move the positions associated with landmark cell synapses until their separation in unrolled phase equals the rest length of the spring, which in turn equals the physical separation between landmark firing fields.

$\phi_A(\tau^-) = \theta_{\text{West}} + k_x L$ . However, upon encountering the east wall, the east border cell pins the attractor phase to  $\theta_{\text{East}}$ .

Before any learning, there is no guarantee that the east border cell pinning phase  $\theta_{\text{East}}$  equals the attractor phase  $\theta_{\text{West}} + k_x L$ , obtained by starting at the west wall and moving to the east wall; sensation and path integration might disagree (Fig. 3B). However, plasticity described in Eq. 10 will act so as to move  $\theta_{\text{East}}$  closer to  $\theta_{\text{West}} + k_x L$ . Then as the animal returns to the left wall at time  $t = 2\tau$ , path integration will retard the attractor phase  $\phi_A(2\tau) = \theta_{\text{East}} - k_x L$ , and an encounter with the west wall leads the west border cell to pin the attractor phase to  $\theta_{\text{West}}$ . Again, there is no guarantee that the west border cell pinning phase  $\theta_{\text{West}}$  agrees with the attractor phase  $\theta_{\text{East}} - k_x L$  obtained by starting at the east wall and traveling to the west wall, but circuit plasticity will change  $\theta_{\text{West}}$  to reduce this discrepancy. Overall, plasticity over multiple cycles of exploration yields the iterative dynamics

$$\theta_{\text{East}} \rightarrow \theta_{\text{West}} + k_x L, \quad \theta_{\text{West}} \rightarrow \theta_{\text{East}} - k_x L.$$

Thus the phase difference  $\theta_{\text{East}} - \theta_{\text{West}}$  between the pinning phases of the two landmark cells will approach the phase advance  $k_x L$  incurred by path integration between the two landmarks. Thus learning can precisely co-register sensation and path integration so that these two information sources yield a consistent map of space (Fig. 3C). In particular, the attractor phase assigned by the composite circuit to any point in the interior of the environment now becomes independent of which direction the animal is traveling, in contrast to the case before learning (compare the assigned interior phases in Fig. 3B versus C).

**Learning as an elastic relaxation between landmarks.** To gain further insight into the learning dynamics, it is useful to interpret the periodic attractor phase  $\phi_A(t)$  as an internal estimate of position through the “unrolled” coordinate variable  $\mathcal{R}^S = \phi_A/k_x$ . Likewise, we can replace the landmark phase  $\theta_i$  with another linear variable  $\mathcal{R}_i^L = \theta_i/k_x$ , denoting the internal representation of the position of landmark  $i$  (Fig. 3D). This enables us to associate physical positions to landmark cells, or more precisely their pinning phases, although these assigned positions are defined only up to shifts of the grid period. Plasticity over the long timescale  $T$  of exploration then yields the following learning dynamics for the physical

positions in unrolled phase for the landmark cells:

$$d\mathcal{R}_{\text{East}}^L/dT = -\mathbf{M}_{\text{WE}} \left[ \mathcal{R}_{\text{East}}^L - (\mathcal{R}_{\text{West}}^L + \Delta\mathcal{R}_{\text{W}\rightarrow\text{E}}^S) \right] \quad (11)$$

$$d\mathcal{R}_{\text{West}}^L/dT = -\mathbf{M}_{\text{WE}} \left[ \mathcal{R}_{\text{West}}^L - (\mathcal{R}_{\text{East}}^L + \Delta\mathcal{R}_{\text{E}\rightarrow\text{W}}^S) \right], \quad (12)$$

where  $\Delta\mathcal{R}_{\text{W}\rightarrow\text{E}}^S = -\Delta\mathcal{R}_{\text{E}\rightarrow\text{W}}^S = L$ .

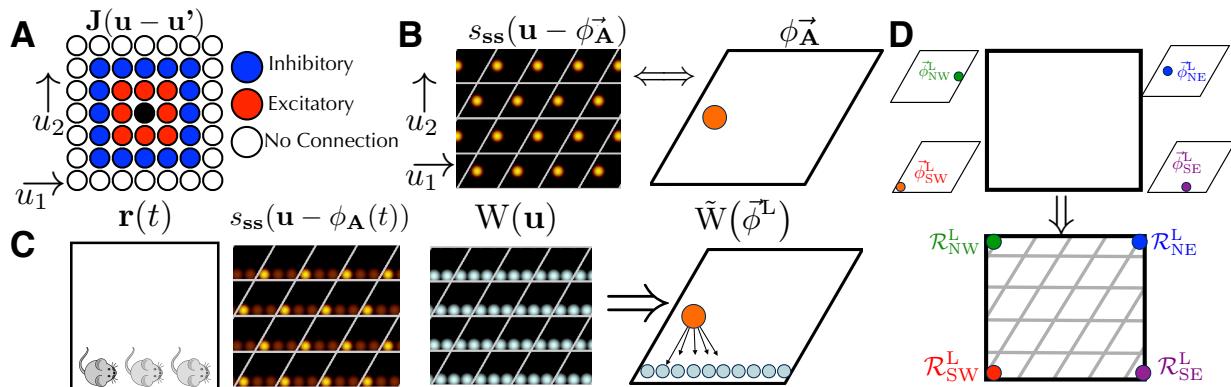
This dynamics for the two landmark cell synapses in unrolled phase is equivalent to that of two particles at physical positions  $\mathcal{R}_{\text{West}}^L$  and  $\mathcal{R}_{\text{East}}^L$ , connected by an overdamped spring with rest length  $L$ , and spring constant  $\mathbf{M}_{\text{WE}}$  which sets the learning rate (Fig. 3E). If the separation  $\mathcal{R}_{\text{East}}^L - \mathcal{R}_{\text{West}}^L$  between the particles is less (greater) than  $L$ , then the spring is compressed (extended) yielding a repulsive (attractive) force between the two particles. Learning stabilizes the two particle positions when their separation equals the spring rest length, so that  $\mathcal{R}_{\text{East}}^L - \mathcal{R}_{\text{West}}^L = L$ . This condition in unrolled phase is equivalent to the fundamental consistency condition for a well defined spatial map, namely that the phase advance due to path integration equals the phase difference between the pinning phases of landmark cells (Fig. 3C). However the utility of the unrolled phase representation lies in revealing a compelling picture for how a spatially consistent map arises from the combined neuronal and synaptic dynamics, through a simple, emergent first order relaxational dynamics of landmark particles connected by damped springs. As we see below, this simple effective particle-spring description of synaptic plasticity in response to spatial exploration generalizes to arbitrary landmarks in arbitrary two dimensional environments.

We note that if the environment has not been fully learned or has been recently deformed, the internal representation of landmarks in unrolled phase will lag behind the true geometry for a time, leading to “boundary-tethered” firing fields seen in (22, 30). Additionally, we have solved the dynamics when the firing fields of the border cells have a finite extent  $L_{\text{wall}}$  and the landmark cells have a finite strength  $\omega$ , and we find the dynamics obeys that of Eq. 11 and Eq. 12, with a different rest length  $\Delta\mathcal{R}_{\text{W}\rightarrow\text{E}}^S = L_{\text{Int}} + 2v \tanh(\omega L_{\text{Wall}}/2v)/\omega$  (See App. D).

## Generalization to 2D Grid Cells

In order to make contact with experiments, we generalize all of the above to two dimensional space. Now grid cells live on a periodic *two-dimensional* neural sheet, where each neuron has position  $\mathbf{u} = (u_1, u_2)$ . The dynamics, analogous to Eq. 1 are:

$$\frac{ds(\mathbf{u})}{dt} = -s(\mathbf{u}) + \mathcal{F} \left[ \iint_{\mathbf{u}'} \mathbf{J}(|\mathbf{u} - \mathbf{u}'|) s(\mathbf{u}') \right]. \quad (13)$$



**Fig. 4.** **A**) A 2D neural sheet with short-range excitation and long-range inhibition, analogous to Fig. 1. Each neuron on the continuous sheet now has coordinates  $\mathbf{u} = (u_1, u_2)$ . **B**) A 2D analogue of a single attractor pattern on the neural sheet, with high firing rates in red (compare to Fig. 1). The set all unique stable attractor patterns is now indexed not by a single phase variable as in 1D, but a 2D phase variable  $\vec{\phi}_A$  ranging over a rhombus, or unit cell. Copies of the unit cell are shown via white lines. **C**) The landmark cell hebbian weights will be a combination of 2D attractor states (Eq. 15). As the animal travels along the south wall, the average firing rates will form a “streak” across the neural sheet. This leads the hebbian weights on the neural sheet to form the same streak; this learned state can be represented as a *distribution* over the periodic rhombus. Analogously, there is a force law, where the state of an attractor network  $\vec{\phi}_A$  will be pulled towards this distribution  $\tilde{W}_i(\vec{\phi}^L)$  (Eq. 14). **D**) Similarly to Fig. 3D, we can unroll the two-dimensional attractor phase into a two-dimensional position variable, thereby associating landmark pinning phases to points in physical space. Given landmarks in all four corners, the landmark pinning phases correspond to different points on the phase rhombus, but through unrolling this rhombus, each can be associated to a physical corner of the environment.

Where  $\mathbf{J}$  includes short-range excitation and long-range inhibition (Fig. 4A). Attractor dynamics on a *two-dimensional* neural sheet can now yield a *two-dimensional* family of stable, or steady state, localized bump activity patterns  $s_{ss}(\mathbf{u} - \vec{\phi}_A)$  with *hexagonal* symmetry (App. E). The attractor state is now a 2D phase  $\vec{\phi}_A$  on the periodic rhombus (Fig. 4B).

Applying the same techniques used to derive Eq. 8, we obtain a 2D analogue to dynamics of the attractor state:

$$\frac{d}{dt}\vec{\phi}_A(t) = \overset{\leftrightarrow}{\mathbf{K}} \cdot \frac{d}{dt}\mathbf{r} + \sum_i \omega_i H_i(\mathbf{r}(t)) \int_{\vec{\phi}^L} \tilde{W}_i(\vec{\phi}^L) \text{Force}_A(\vec{\phi}^L - \vec{\phi}_A). \quad (14)$$

Here we replaced  $k_x$  in 1D with  $\overset{\leftrightarrow}{\mathbf{K}}$ , a matrix that translates 2D animal velocity into phase advance in the 2D attractor network;  $\overset{\leftrightarrow}{\mathbf{K}}$  determines *both* grid spacing *and* orientation. The analog of learning dynamics in (Eq. 7) is:

$$\frac{d\tilde{W}_i(\vec{\phi}^L)}{dT} = \Pr(\vec{\phi}_A(t) = \vec{\phi}^L | i \text{ Firing}) - \tilde{W}_i(\vec{\phi}^L), \quad (15)$$

where  $\tilde{W}_i(\vec{\phi}^L)$  is now a distribution over the periodic rhombus (Fig. 4C).

In an analogous manner, we may make a small-angle approximation to replace the attractor phase  $\phi_A(t)$  with a *two-dimensional* attractor coordinate variable  $\mathcal{R}^S(t)$ , reflecting an internal estimate of instantaneous position in physical space, and we replace the landmark phase  $\tilde{W}_i(\vec{\phi}^L)$  with another 2D attractor coordinate variable  $\mathcal{R}_i^L$

reflecting the internal estimate of landmark position (Fig. 4D). This yields two-dimensional dynamics for position self-estimates and landmark position, given in analogy to Eqs. 9, 10 by:

$$d\mathcal{R}^S/dt = d\mathbf{r}/dt + \sum \omega_i H_i(\mathbf{r}(t)) \cdot (\mathcal{R}_i^L - \mathcal{R}^S), \quad (16)$$

$$d\mathcal{R}_i^L/dT = \langle \mathcal{R}^S(t) | \text{Cell } i \text{ Firing} \rangle - \mathcal{R}_i^L. \quad (17)$$

## Spatial consistency through emergent elasticity

We showed in Eq. 11 and Eq. 12, and in Fig. 3DE that the emergence of spatial consistency between path integration and landmarks through Hebbian learning dynamics, during exploration of a simple 1D environment, could be understood as the outcome of an elastic relaxation process between landmark cell synapses, viewed as particles in physical space connected by damped springs. Remarkably, this result generalizes far beyond this simple environment. As long as the exploration dynamics are time-reversible<sup>1</sup>, the learning dynamics of *any* set of landmark cells in *any* geometry yields this particle-spring interpretation:

$$d\mathcal{R}_i^L/dT = - \sum_j \mathbf{M}_{ij} \left( \mathcal{R}_i^L - [\mathcal{R}_j^L + \Delta\mathcal{R}_{j \rightarrow i}] \right). \quad (18)$$

The spring constant  $\mathbf{M}_{ij}$  is related to the frequency with which the animal moves between each pair of landmark

<sup>1</sup>Time-reversible means that for any  $\mathbf{r}(t)$ , the reverse path  $\mathbf{r}(-t)$  is equally likely

firing fields  $i, j$ , while the rest displacement  $\Delta\mathcal{R}_{j \rightarrow i}^S$  is the average change in unrolled attractor phase as the animal moves from firing field  $j$  to field  $i$ , roughly related to the distance between the landmark firing fields. Precise expressions for the spring constants and rest lengths are derived from the statistics of exploration in (App. F). Overall, this elastic relaxation process converges towards an internal map where all pairs of landmark cell synapses, viewed as particles in unrolled phase, or physical space, become separated by the physical distance between their firing fields. This convergence ensures a consistent internal environmental map of external space in which velocity based path integration of attractor phase starting at the pinning phase of landmark  $i$  and ending at landmark  $j$  will yield an integrated phase consistent with the pinning phase of landmark  $j$  itself. This relaxation dynamics explains path-dependent shifts in firing patterns observed in recently deformed environments (22). Also, the experimental observation that in multi-compartment environments, consistent maps *within* compartments form before consistent maps *between* compartments are also explained (10) by this relaxation dynamics. In essence, the longest-lived learning mode of the relaxation dynamics corresponds to differences in maps *between* compartments.

Furthermore, as we explain in the next three sections, these relaxation dynamics yields several novel experimental predictions: (1) systematic path-dependent shifts in *fully learned* 2D environments, (2) mechanical deformations in complex environments, and (3) the novel prediction of creation of topological defects in grid cell firing patterns through specific environmental manipulations.

## Path-dependence in 2D environments

We saw above that exploration in a simple 1D geometry lead to a consistent internal map in which the attractor network phase was mapped onto the current physical position alone, independent of path history (Fig. 3C). This consistency arises through the elastic relaxation process in Eq. 11 and Eq. 12, which makes the distance between the landmark cells in unrolled phase  $\mathcal{R}_{\text{East}}^L - \mathcal{R}_{\text{West}}^L$  equal to the physical distance between their firing fields  $L$ , just like two particles connected by a spring with rest length  $L$  (Fig. 3E). This situation will generalize to two dimensions if there are only two landmarks, namely a west and east border cell (Fig. 5A1). However, it becomes more complex with the addition of a third landmark cell, for example a south border cell (Fig. 5A2).

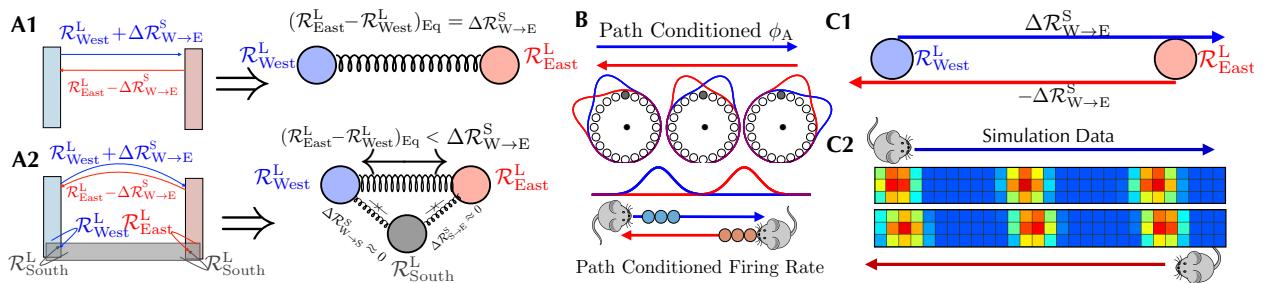
In this case, east and west landmark particles will be connected by a spring of rest length  $\Delta\mathcal{R}_{\text{E} \rightarrow \text{W}}^S = L$ , as before, but they will each also be connected to the south

landmark particle with springs. Intuitively, as the mouse travels from the east or west walls to the south walls, the landmark pinning phases of each of these three border cells will be attracted towards each other<sup>2</sup>. The combined three particle elastic system will settle into an equilibrium configuration in which the difference in unrolled phase between east and west landmarks will be *less* than the physical separation  $L$ , or equivalently the rest length  $\Delta\mathcal{R}_{\text{E} \rightarrow \text{W}}^S$  of the spring connecting them. This in turn implies that the attractor phase assigned to any physical position in the interior will be relatively phase advanced (retarded) if the mouse is on a trajectory leaving the west (east) wall. This path dependence in the attractor phase is entirely analogous to that seen in Fig. 3B. However, the reason is completely different. In Fig. 3B, the landmark particles are not separated by the rest length of the spring connecting them because the environment is not fully learned and so the particles are out of equilibrium, whereas in Fig. 5A2, the particles are not separated by the rest length, even in a *fully learned* environment, because additional springs from the south landmark create excess compression.

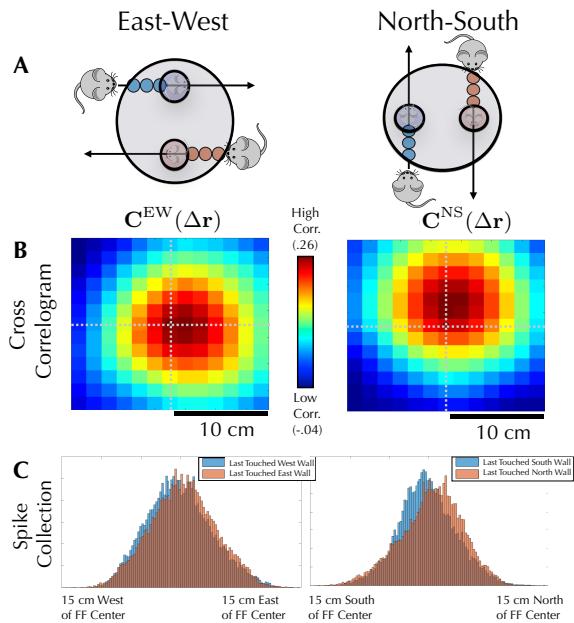
This theory makes a striking experimentally testable prediction, namely that even in a *fully learned* 2D environment, grid cell firing fields, when computed on subsets of mouse trajectories conditioned on leaving a particular border, will be shifted *towards* that border (Fig. 5B). This shift occurs because at any given position, the attractor phase depends on the most recently encountered landmark. In particular, on a west to east (east to west) trajectory, the attractor phase will be advanced (retarded) relative to a east to west (west to east) trajectory. Thus on a west to east trajectory, the advanced phase will cause grid cells to fire earlier, yielding west shifted grid cell firing fields as a function of position. Similarly on an east to west trajectory, grid fields will be east shifted. In summary, the theory predicts grid cell firing patterns conditioned on trajectories leaving the west (east) border will be shifted west (east). While we have derived this prediction qualitatively using the conceptual mass-spring picture in Fig. 5A2, we confirm this intuition through direct numerical simulations of the full circuit dynamics in Eq. 14 and Eq. 15 (Fig. 5C2). Under reasonable parameters, our simulations can yield path-dependent shifts of up to  $\sim 2$  cm towards whichever wall the animal last touched (App. G).

We searched for such subtle shifts in a population of 143 grid cells from 14 different mice that had been exploring a familiar, well-learned, 1-meter open field (App. H), using

<sup>2</sup>More complex, non-overlapping distributions yield the same deformations.



**Fig. 5. A1)** For two landmark cells, the rest length  $\Delta\mathcal{R}_{W \rightarrow E}^S$  of the spring connecting them equals the physical width  $L$  of the environment, and so the two landmark particles learn unrolled pinning phases  $\mathcal{R}_{East}^L$  and  $\mathcal{R}_{West}^L$  obeying the spatial consistency condition  $(\mathcal{R}_{East}^L - \mathcal{R}_{West}^L)_{Eq} = \Delta\mathcal{R}_{W \rightarrow E}^S = L$  as in Fig. 3C. **A2)** The addition of a southern landmark cell will cause a pinning effect which pulls  $\mathcal{R}_{West}^L$ ,  $\mathcal{R}_{East}^L$  closer together. The animal can travel from the east and west landmark field to the southern landmark field with little path integration at all, yielding  $\Delta\mathcal{R}_{W \rightarrow S}^S \approx 0$ ,  $\Delta\mathcal{R}_{S \rightarrow E}^S \approx 0$ . **B)** If the attractor phase is advanced on a west to east trajectory (blue) relative to an east to west trajectory (red), then any particular grid cell (in this case the shaded grey cell) will fire earlier (later) on west-to-east (east-to-west) trajectory. Thus grid fields computed from trajectories leaving the west (east) border will shift west (east). **C1)** When landmark pinning phases are pulled together closer than the path integration distance between them, then the attractor phase will shift away from whichever wall the animal last encountered. Therefore it will phase advance on west-to-east trajectories relative to east-to-west trajectories, as in Fig. 3B and Fig. 5B. **C2)** Thus simulations of Eq. 14 and Eq. 15 lead to grid cell firing patterns shifted *towards* whichever wall the animal last encountered.



**Fig. 6. A)** Our theory predicts that grid cell firing patterns will be shifted towards whichever wall the animal last encountered, even in a fully learned environment. **B)** On the left (right) this shift is detected by computing the cross correlation between west (south) conditioned firing fields, shifted by a spatial offset  $\Delta r$ , and the unshifted east (north) conditioned firing field. The cross correlation peaks when the spatial shift is positive in the x-direction (positive in the y-direction), as predicted by theory. **C)** This effect can also be seen by comparing histograms of spike positions around firing field centers for different path conditions.

two separate analyses, based on cross-correlations and spike shifts with respect to field centers.

**Cross-Correlations** One method for detecting a systematic firing field shift across many grid fields is to cross-correlate firing rate maps conditioned on trajectories leaving two different borders (App. A.1). For ex-

ample, for each cell, we can ask how much and in what direction we must shift its west border conditioned firing field to match, or correlate as much as possible with, the same cell's east conditioned firing field. In particular, for each cell, we can compute the correlation coefficient between a spatially shifted west conditioned field and an unshifted east conditioned field, and plot the average correlation coefficient as a function of this spatial shift. The theory predicts that we will have to shift the west conditioned firing field eastwards to match the east conditioned firing field (Fig. 6A). This prediction is confirmed by a peak in the cross-correlation as a function of spatial shift when the spatial shift is positive, or directed east (Fig. 6B).

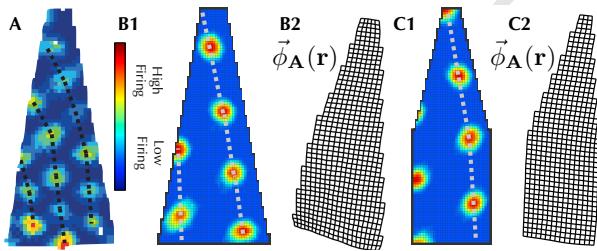
A similar logic holds for north and south; in order to maximally correlate the south conditioned field to the north conditioned field, the theory predicts we will need to shift the south conditioned map north. This requisite shift is seen in the data in Fig. 6B, which reflects the cross correlation between a shifted south field and an unshifted north field, averaged across all cells. The maximal correlation is achieved when the south fields are shifted north. Overall, this analysis shows that grid patterns are shifted towards the most recently encountered wall, both for the NS walls ( $3 \text{ cm}$ ,  $P = 1.5 \cdot 10^{-5}$ , Binomial Test,  $P = 1.5 \cdot 10^{-5}$ , Sign-Flip Test) and the EW Walls ( $1.5 \text{ cm}$ ,  $P = 10^{-7}$ , Binomial Test,  $P = 10^{-7}$ , Sign-Flip Test), matching the sign and magnitude seen in simulations.

**Firing Field Centers** These results can be corroborated by computing shifts in spikes relative to firing field centers, when conditioning spikes on the path history (App. A.2). For each firing field center, we calculate the average spike position within that firing field conditioned on the animal having last touched a particular wall. For each

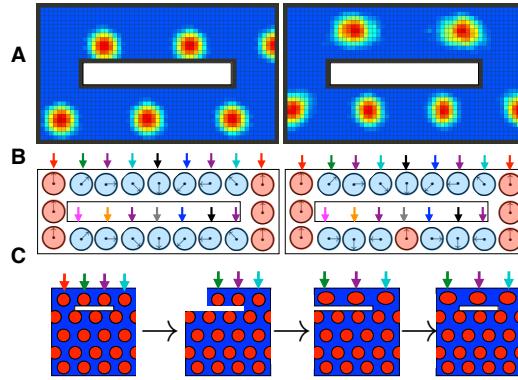
cell, we calculate the average shift across all firing fields, and examine how the shifts depend on which wall the animal last touched. Again, the patterns are shifted towards whichever wall the animal last touched (Fig. 6C) for both the NS walls (.5 cm,  $P = 10^{-5}$  Binomial Test,  $P = 10^{-5}$  Sign-Flip Test) and the EW Walls (.5 cm,  $P = 3 \cdot 10^{-4}$  Binomial Test,  $P = 2 \cdot 10^{-2}$  Sign-Flip Test). The discrepancy in the estimated magnitude of the shift between the methods of analysis is likely due to poorly defined firing fields; a method based on firing field centers will give a lower signal-to-noise ratio, and thus a lower shift magnitude, than the cross-correlogram method.

## Mechanical deformations in complex environments

Another experimental observation that can be reproduced by our theory is the distortion (11) of grid cell patterns seen in an irregular environment (Fig. 7A). Landmark cells with firing fields distributed across an entire wall will pull the attractor phase to its associated landmark pinning phase *regardless* of where along the wall the animal is. The presence of a diagonal wall then causes the average attractor phase as a function of position to curve towards the wall, yielding spatial grid cell patterns that curve *away* from the wall (Fig. 7B, C). Previous theoretical accounts of this grid cell deformation have relied on purely phenomenological models that treated individual grid cell firing fields as particles with mostly repulsive interactions (15), without a clear mechanistic basis underlying this interaction. Here we provide, to our knowledge for the first time, a model with a clear mechanistic basis for such deformations, grounded in the interaction between attractor based path integration and landmark cells with plastic synapses. Such dynamics yields an emergent elasticity where the particles are landmark cell synapses rather than individual firing field centers.



**Fig. 7.** **A)** Experimental data of grid cell firing patterns deformed, curving away from a wall in an irregular geometry. **B1)** A full simulation of Eq. 14, Eq. 15 also yields grid firing patterns bent away from the wall. **B2)** Visualization of the average attractor state as a function of position  $\vec{\phi}_A(r)$  (periodicity removed for visualization purposes). The reversal between the bending of the internal attractor phase and the bending of firing rate maps is similar to the reversal seen in Fig. 5B. **C1), C2)** Same as B1), B2), but for a slightly different geometry.



**Fig. 8.** **A)** Two steady state grid cell patterns emerging from the same cue-rich environment. In the first firing pattern, the combination of landmark pinning and path integration yields a phase advance of four firing fields in traveling from west to east along either corridor. The second pattern has a topological defect; traveling from the west to east through the *north* corridor yields a phase increase of  $\sim 1.5$  firing fields; traveling east to west through the *south* corridor yields a phase decrease of  $\sim 2.5$  firing fields. This second pattern is stable nonetheless. **B)** Schematic of 1D underlying attractor state as a function of space. The two patterns in (A) correspond to two different landmark pinning phase patterns learned by the many landmarks. Both landmark pinning patterns are stable under Eq. 14, Eq. 15. In the first pattern, the combination of landmark pinning and path integration yields the same phase advance in both the north and south corridors. The second pattern has a topological defect; the phase advance in the north corridor *one full rotation less* than the phase advance through the south corridor. This is possible because many landmark cues (colored arrows) can yield many landmark cells with multiple stable synaptic configurations, or pinning phases under Eq. 14, Eq. 15. **C)** Schematic of proposed “deformation schedule” that could yield a topological defect in grid cell firing patterns. By separating/truncating the northern corridor, stretching it (along with spatial cues, denoted by colored arrows), and reconnecting it, it may be possible introduce one of these defects. Even though the initial geometry is identical to the final geometry, the deformation schedule has lead to a firing pattern which is three fields wide in the north and four fields wide in the south.

## Topological defects in grid cells: a prediction

While the dynamics of the linearized Eq. 17 will always flow to the same relative landmark representations  $\mathcal{R}_i^L$ , this is not the case for the full dynamics of Eq. 14, Eq. 15, which can learn multiple different stable landmark cell synaptic configurations. One striking example of this is the ability of the learning dynamics to generate “topological defects”, where the number of firing fields traversed is not the same for two different paths (Fig. 8A, B and App. G). An environmental geometry capable of supporting these defects will yield a set of firing patterns that depends not only on the final geometry, but also on the *history* of how this geometry was created (Fig. 8C).

## Discussion

Overall, we have provided a theoretical framework for exploring how sensory cues and path integration may

work together to create a consistent internal representation of space. Our framework is grounded in biologically plausible mechanisms involving attractor based path integration of velocity and Hebbian plasticity of landmark cells. Moreover, systematic model reduction of this combined neural and synaptic dynamics yields a simple and intuitive emergent elasticity model in which landmark cell synapses act like particles sitting in physical space connected by damped springs whose rest length is equal to the physical distance between landmark firing fields. This simple emergent elasticity model not only provides a conceptual explanation of how neuronal dynamics and synaptic plasticity can conspire to self-organize a consistent map of space in which sensory cues and path-integration are in register, but also provides novel predictions involving small shifts in firing fields even in fully learned environments, the possibility of topological defects in grid cells, and the mechanical deformation of grid cells in response to irregular borders.

This work opens up many interesting avenues for future research. For example, further explorations of the nonlinear regime of our combined circuit dynamics may yield interesting experimental signatures that distinguish different modes of interactions between attractor networks, path integrators and landmark cells. Incorporating heterogeneity of neural representations observed in MEC (31) into our framework is another intriguing avenue. Also, as the reliability of sensory and velocity cues change, it is interesting to ask what higher order mechanisms may exist to differentially regulate the effect of landmarks and velocity on the internal representation of space. More generally, our theory provides a unified framework for understanding how systematic variations in environmental geometry and the statistics of environmental exploration interact to precisely sculpt neural representations of space.

**Acknowledgements** LMG is a New York Stem Cell Foundation - Robertson Investigator. This work was supported by funding from The New York Stem Cell Foundation, James S McDonnell Foundation, Whitehall Foundation, NIMH MH106475 and a Klingenstein-Simons Fellowship awarded to LMG, funding from the Simons Foundation awarded to LMG and SG. SO was supported by the Karel Urbanek Postdoctoral Fellowship in Applied Physics. KH is supported by a Stanford SIGF. We thank Daniel Fisher for helpful discussions.

1. Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic robotics*. MIT press, 2005.
2. Torkel Hafting, Marianne Fyhn, Sturla Molden, May-Britt B Moser, and Edvard I Moser. Microstructure of a spatial map in the entorhinal cortex. *Nature*, 436(7052):801–6, 2005. ISSN 0028-0836. doi: 10.1038/nature03721.
3. Edvard I. Moser, Yasser Roudi, Menno P. Witter, Clifford Kentros, Tobias Bonhoeff-

fer, and May-Britt Moser. Grid cells and cortical representation. *Nature Reviews Neuroscience*, 15, 06 2014.

4. Trygve Solstad, Charlotte N Boccara, Emilio Kropff, May-Britt B Moser, and Edvard I Moser. Representation of geometric borders in the entorhinal cortex. *Science*, 322 (5909):1865–8, 2008. ISSN 0036-8075. doi: 10.1126/science.1166466.
5. Emilio Kropff, James E Carmichael, May-Britt Moser, and Edvard I Moser. Speed cells in the medial entorhinal cortex. *Nature*, 523(7561):419–424, 2015.
6. Francesca Sargolini, Marianne Fyhn, Torkel Hafting, Bruce L McNaughton, Menno P Witter, May-Britt Moser, and Edvard I Moser. Conjunctive representation of position, direction, and velocity in entorhinal cortex. *Science*, 312(5774):758–762, 2006.
7. Tor Stensola, Hanne Stensola, May-Britt Moser, and Edvard I. Moser. Shearing-induced asymmetry in entorhinal grid cells. *Nature*, 518, 02 2015.
8. Caswell Barry, Robin Hayman, Neil Burgess, and Kathryn J Jeffery. Experience-dependent rescaling of entorhinal grids. *Nature neuroscience*, 10(6):682–684, 2007. ISSN 1097-6256. doi: 10.1038/nn1905.
9. Caswell Barry, Lin Ginzberg, O’Keefe, John, and Neil Burgess. Grid cell firing patterns signal environmental novelty by expansion. *Proceedings of the ...*, 109(43):17687–17692, 2012. ISSN 0027-8424. doi: 10.1073/pnas.1209918109.
10. Francis Carpenter, Daniel Manson, Kate Jeffery, Neil Burgess, and Caswell Barry. Grid cells form a global representation of connected environments. *Current Biology*, 25(9):1176–1182, 2014/04/04 2015. doi: 10.1016/j.cub.2015.02.037.
11. Julija Krupic, Marius Bauza, Stephen Burton, Caswell Barry, and O’Keefe, John. Grid cell symmetry is zenvironmental geometry. *Nature*, 518(7538), 2015. ISSN 0028-0836. doi: 10.1038/nature14153.
12. Julija Krupic, Marius Bauza, Stephen Burton, and John O’Keefe. Local transformations of the hippocampal cognitive map. *Science*, 359(6380):1143–1146, 2018. ISSN 0036-8075. doi: 10.1126/science.aao4960.
13. Julija Krupic. Brain crystals. *Science*, 350(6256):47–48, 2015. ISSN 0036-8075. doi: 10.1126/science.aad3002.
14. Julija Krupic, Marius Bauza, Stephen Burton, and John O’Keefe. Framing the grid: effect of boundaries on grid cells and navigation. *The Journal of physiology*, 594(22):6489–6499, 2016.
15. Julija Krupic, Marius Bauza, Stephen Burton, Colin Lever, and O’Keefe, John. How environment geometry affects grid cell symmetry and what we can learn from it. ... of the Royal ... , 369(1635):20130188, 2014. ISSN 0962-8436. doi: 10.1098/rstb.2013.0188.
16. Bruce L McNaughton, Francesco P Battaglia, Ole Jensen, Edvard I Moser, and May-Britt Moser. Path integration and the neural basis of the cognitive map. *Nature Reviews Neuroscience*, 7(8):663–678, 2006.
17. Taifan Evans, Andrej Bicanski, Daniel Bush, and Neil Burgess. How environment and self-motion combine in neural representations of space. *The Journal of Physiology*, 594(22):6535–6546, 2016. doi: 10.1113/JP270666.
18. Florian Raudies, James R Hinman, and Michael E Hasselmo. Modelling effects on grid cells of sensory input during self-motion. *The Journal of physiology*, 594(22):6513–6526, 2016.
19. Marianne Fyhn, Torkel Hafting, Alessandro Treves, May-Britt B Moser, and Edvard I Moser. Hippocampal remapping and grid realignment in entorhinal cortex. *Nature*, 446(7132):190–4, 2007. ISSN 0028-0836. doi: 10.1038/nature05601.
20. Kiah Hardcastle, Surya Ganguli, and Lisa M Giocomo. Environmental boundaries as an error correction mechanism for grid cells. *Neuron*, 86(3):827–839, 2015.
21. Lisa M Giocomo. Environmental boundaries as a mechanism for correcting and anchoring spatial maps. *The Journal of Physiology*, 594(22):6501–6511, 11 2016. doi: 10.1113/JP270624.
22. Alex T. Keinath, Russell A. Epstein, and Vijay Balasubramanian. Environmental deformations dynamically shift the spatial metric of the brain. *bioRxiv*, 2017. doi: 10.1101/174367.
23. Samuel Ocko, Kiah Hardcastle, Lisa Giocomo, and Surya Ganguli. Evidence for optimal bayesian cue combination of landmarks and velocity in the entorhinal cortex. *Cosyne*, 2017.
24. Michael J Milford, Gordon F Wyeth, and David Prasser. Ratslam: a hippocampal model for simultaneous localization and mapping. In *Robotics and Automation, 2004. Proceedings. ICRA’04. 2004 IEEE International Conference on*, volume 1, pages 403–408. IEEE, 2004.
25. M Mulas, N Waniek, and J Conradt. Hebbian plasticity realigns grid cell activity with external sensory cues in continuous attractor models. *Frontiers in computational ...*, 2016.
26. Gustavo Deco, Viktor K. Jirsa, Peter A. Robinson, Michael Breakspear, and Karl Friston. The dynamic brain: From spiking neurons to neural masses and cortical fields. *PLOS Computational Biology*, 4(8):1–35, 08 2008. doi: 10.1371/journal.pcbi.1000092.
27. Alexei Samsonovich and Bruce L. McNaughton. Path integration and cognitive mapping in a continuous attractor neural network model. *Journal of Neuroscience*, 17 (15):5900–5920, 1997. ISSN 0270-6474. doi: 10.1523/JNEUROSCI.17-15-05900. 1997.
28. Yoram Burak and Ila R Fiete. Accurate path integration in continuous attractor network models of grid cells. *PLoS Comput Biol*, 5(2):e1000291, 2009.
29. Guifen Chen, Daniel Manson, Francesca Cacucci, and Thomas Joseph Wills. Ab-

- sence of visual input results in the disruption of grid cell firing in the mouse. *Current Biology*, 26(17):2335–2342, 2016.
- 30. KM Gothard, WE Skaggs, and Bruce L. Dynamics of mismatch correction in the hippocampal ensemble code for space: interaction between path integration and environmental cues. *J. Neurosci.*, 16(24):8027–40, 1996. ISSN 0270-6474.
  - 31. Kiah Hardcastle, Niru Maheswaranathan, Surya Ganguli, and Lisa M. Giocomo. A multiplexed, heterogeneous, and adaptive code for navigation in medial entorhinal cortex. *Neuron*, 94(2):375 – 387.e7, 2017. ISSN 0896-6273. doi: <https://doi.org/10.1016/j.neuron.2017.03.025>.

DRAFT

## Supplementary Note A: Reducing perturbation to an effective force law in the one-variable “ring” representation

Consider a steady-state firing rate pattern  $s_{SS}(u - \phi_A)$  under some translation-invariant neural dynamics function  $\text{Dyn}$  that is given some perturbation centered at  $\phi^P$  having the same periodicity as  $s_{SS}$ :

$$\frac{ds(u)}{dt} = \text{Dyn}[s(u)] + \epsilon \text{Pert}(u - \phi^P).$$

In order to understand the effect of this perturbation, we need to understand the Jacobian matrix around the point  $s_{SS}(u - \phi_A)$ :

$$\left. \frac{d\text{Dyn}}{ds} \right|_{s=s_{SS}(u-\phi_A)+\Delta s} \approx \mathbf{Jac}^{\phi_A} \cdot \Delta s$$

**A. Modes of the Jacobian.** Because  $s_{SS}(u - \phi_A)$  is a stable one-dimensional family of solutions of  $\text{Dyn}$ ,  $\mathbf{Jac}^{\phi_A}$  must be a negative semidefinite matrix. Because

$$\text{Dyn}[s_{SS}(u - \phi_A)] = 0 \text{ For all } \phi_A.$$

there is a single-zero eigenvector<sup>3</sup>, the sliding mode  $\frac{ds_{SS}(u-\phi_A)}{du}$ :

$$\frac{d\text{Dyn}[s_{SS}(u - \phi_A)]}{d\phi_A} = \left. \frac{d\text{Dyn}[s(u)]}{d\left[\frac{ds_{SS}(u-\phi_A)}{du}\right]} \right|_{s=s_{SS}(u-\phi_A)} = \mathbf{Jac}^{\phi_A} \cdot \left[ \frac{ds_{SS}(u - \phi_A)}{du} \right] = 0.$$

**B. Effect of small perturbations.** When an external perturbation is small and  $\mathbf{Jac}^{\phi_A}$  is symmetric, e.g.  $\frac{d\text{Dyn}[s](u')}{ds(u')} = \frac{d\text{Dyn}[s](u)}{ds(u')}$ , the effective perturbation will be the projection of the actual perturbation onto the sliding mode.

$$\begin{aligned} \frac{ds}{dt} &\approx \epsilon \underbrace{\left[ \int \frac{ds_{SS}(u - \phi_A)}{du} \text{Pert}(u - \phi^P) \right]}_{\text{Projection onto Sliding Mode}} \cdot \underbrace{\frac{ds_{SS}(u - \phi_A)}{du}}_{\text{Sliding Mode}} \\ &= \epsilon \underbrace{\left[ \int \frac{ds_{SS}(u)}{du} \text{Pert}(u - [\phi^P - \phi_A]) \right]}_{-\text{Force}_A(\phi^P - \phi_A) \text{ (Definition)}} \underbrace{\frac{ds_{SS}(u - \phi_A)}{du}}_{\text{Sliding mode}} \\ &= -\epsilon \text{Force}_A(\phi^P - \phi_A) \frac{ds_{SS}(u - \phi_A)}{du}. \end{aligned}$$

We can translate these dynamics into the reduced  $\phi$  representation:

$$\frac{d\phi_A}{dt} = \epsilon \text{Force}_A(\phi^P - \phi_A), \quad (19)$$

Eq. 19 can be verified:

$$\frac{ds_{SS}(u - \phi_A)}{dt} = \frac{ds_{SS}(u - \phi_A)}{d\phi_A} \frac{d\phi_A}{dt} = -\frac{ds_{SS}(u - \phi_A)}{du} \frac{d\phi_A}{dt} = -\underbrace{\epsilon \text{Force}_A(\phi^P - \phi_A)}_{d\phi_A/dt} \frac{ds_{SS}(u - \phi_A)}{du} \quad (20)$$

<sup>3</sup>We can show this by contradiction; if  $\mathbf{Jac}^{\phi_A}$  had any other non-negative modes, the family of steady states would be larger

**B.1. Form of the landmark cell force function.** When the perturbation takes the form of input from Hebbian landmark cells, the perturbation function is simply the attractor bump pattern  $s_{\text{SS}}(u - \phi^L)$ . Therefore, defining  $\Delta\phi = \phi^L - \phi_A$ ,

$$\text{Force}_A(\Delta\phi) = - \int_u \frac{ds_{\text{SS}}(u)}{du} s_{\text{SS}}(u - \Delta\phi) = - \frac{d}{d\Delta\phi} \left[ \int_u s_{\text{SS}}(u) s_{\text{SS}}(u - \Delta\phi) \right]$$

Therefore, the force function is simply the negative derivative of the spatial autocorrelation function of the bump pattern. Because the spatial autocorrelation is even and maximized at  $\Delta\phi = 0$ , minimized at  $\Delta\phi = \pi$ , the force function will be odd, with positive(negative) values for positive(negative)  $\Delta\phi$ . As long as the bump size is not much smaller than the bump spacing, the autocorrelation will decrease gradually between  $\Delta\phi = 0, \Delta\phi = \pi$ , leading to a long range force function which only approaches zero at, and far from, the origin. This behavior is qualitatively matched by  $\text{Force}_A(\Delta\phi) = \sin(\Delta\phi)$ . For simplicity, we define the magnitude of  $\text{Force}_A(\Delta\phi)$  to give it a slope of 1 at  $\Delta\phi = 0$ ; all strength information can be contained in  $\omega$ .

**B.2. Dynamics with non-symmetric Jacobians.** When  $\text{Jac}^{\phi_A}$  is non-symmetric, we may use the same techniques as before, except now we must use a non-orthogonal projection onto the sliding mode:

$$\frac{ds}{dt} \approx \epsilon \underbrace{\left[ \int_u \mathbf{v}_{\text{proj}}(u - \phi_A) \text{Pert}(u - \phi^P) \right]}_{\text{Non-orthogonal projection onto sliding mode}} \cdot \underbrace{\frac{ds_{\text{SS}}(u - \phi_A)}{du}}_{\text{Sliding Mode}},$$

where  $\mathbf{v}_{\text{proj}}(u - \phi_A)$  can, in principle, be solved through diagonalization of the Jacobian  $\text{Jac}^{\phi_A}$ .

## Supplementary Note B: Proof of recovery of exact path integration

By coupling the attractor network to conjunctive position and velocity-tuned cells that east (west) movement-selective cells form feedforward synapses into the attractor network that are shifted in the positive (negative)  $u$  (28) direction, we impose a velocity-dependent perturbation on the network. This yields,

$$\begin{aligned} ds/dt &= \text{Dyn}[s(u)] + v_{\text{East}} \epsilon \text{Pert}(u - [\phi_A + \Delta\phi_{\text{PI}}]) \\ &\quad + v_{\text{West}} \epsilon \text{Pert}(u - [\phi_A - \Delta\phi_{\text{PI}}]), \end{aligned}$$

where  $v_{\text{East}}, v_{\text{West}}$  are the east and west velocities of the animal. Then model reduction via Eq. 2 yields

$$\begin{aligned} d\phi_A/dt &= v_{\text{East}} \epsilon \text{Force}_A(\Delta\phi_{\text{PI}}) + v_{\text{West}} \epsilon \text{Force}_A(-\Delta\phi_{\text{PI}}) = \\ &\quad \underbrace{[v_{\text{East}} - v_{\text{West}}]}_{v_x = dr_x/dt} \underbrace{\epsilon \text{Force}_A(\Delta\phi_{\text{PI}})}_{k_x (\text{Definition})} = v_x \cdot k_x. \end{aligned}$$

Here  $k_x$  is a constant of proportionality that relates animal velocity to the rate of phase advance in the attractor network.

## Supplementary Note C: Lemmas about Landmark Cells

**A. Verifying the Hebbian learning rule in the attractor basis.** We can verify that in the attractor basis,

$$W_i(u) = \int_{\phi^L} \tilde{W}_i(\phi^L) s_{\text{SS}}(u - \phi^L) \tag{21}$$

the learning rule:

$$\frac{d\tilde{W}_i(\phi^L)}{dT} = \text{Pr}(\phi^L | i \text{ Firing}) - \tilde{W}_i(\phi^L) \tag{22}$$

gives us the learning rule in the neural basis:

$$\frac{d\mathbf{W}_i(u)}{dT} = \langle s(u)|i \text{ Firing}\rangle - \mathbf{W}_i(u) = \int_{\phi} s_{\mathbf{SS}}(u - \phi^L) \Pr(\phi^L|i \text{ Firing}) - \mathbf{W}_i(u) \quad (23)$$

by inspection:

$$\begin{aligned} \frac{d\mathbf{W}_i(u)}{dT} &\stackrel{\text{Basis Switch}}{=} \frac{d}{dT} \left[ \int_{\phi} \tilde{\mathbf{W}}_i(\phi^L) s_{\mathbf{SS}}(u - \phi) \right] = \int_{\phi^L} \frac{d\tilde{\mathbf{W}}_i(\phi^L)}{dT} s_{\mathbf{SS}}(u - \phi^L) \\ &\stackrel{\text{Eq.22}}{=} \int_{\phi^L} [\Pr(\phi^L|i \text{ Firing}) - \tilde{\mathbf{W}}_i(\phi^L)] s_{\mathbf{SS}}(u - \phi^L) = \int_{\phi^L} \Pr(\phi^L|i \text{ Firing}) s_{\mathbf{SS}}(u - \phi) - \underbrace{\int_{\phi} \tilde{\mathbf{W}}_i(\phi^L) s_{\mathbf{SS}}(u - \phi^L)}_{\mathbf{W}(u)} \\ &= \underbrace{\int_{\phi} s_{\mathbf{SS}}(u - \phi^L) \Pr(\phi^L|i \text{ Firing}) - \mathbf{W}_i(u)}_{\text{Eq.23}} \end{aligned}$$

## B. Lemmas about linear approximations .

**B.1. Linear Approximation of Forcing Rule .** We can show how linearizing the force law into a simple “spring constant” allows us to represent the landmark state with a single number  $\theta = \int_{\phi^L} \tilde{\mathbf{W}}_i(\phi^L) \phi^L$ :

$$\begin{aligned} \int_{\phi^L} \tilde{\mathbf{W}}_i(\phi^L) \text{Force}_A[\phi^L - \phi_A] &\approx \int_{\phi^L} \tilde{\mathbf{W}}_i(\phi^L) \cdot [\phi^L - \phi_A] \underbrace{\left[ \frac{d}{d\phi} \text{Force}_A(\phi) \Big|_{\phi=0} \right]}_{\text{(Definition of )}} \\ &= \left( \underbrace{\int_{\phi^L} \tilde{\mathbf{W}}_i(\phi^L) \phi^L}_{\theta_i} - \underbrace{\int_{\phi^L} \tilde{\mathbf{W}}_i(\phi^L) \phi_A}_{\phi_A} \right) = \cdot (\theta_i - \phi_A). \end{aligned}$$

**B.2. Proof of single-variable representation of landmarks .** Combining the attractor-basis dynamics for hebbian learning and the linear approximation:

$$\frac{d\tilde{\mathbf{W}}_i(\phi^L)}{dT} = \Pr(\phi^L|i \text{ Firing}) - \tilde{\mathbf{W}}_i(\phi^L), \quad \theta = \int_{\phi^L} \tilde{\mathbf{W}}_i(\phi^L) \phi^L,$$

we get:

$$\begin{aligned} \frac{d\theta_i}{dT} &= \int_{\phi^L} \frac{d\tilde{\mathbf{W}}_i(\phi^L)}{dT} \phi^L = \int_{\phi^L} [\Pr(\phi^L|i \text{ Firing}) - \tilde{\mathbf{W}}_i(\phi^L) \phi^L] \\ &= \underbrace{\int_{\phi^L} \Pr(\phi^L|i \text{ Firing}) \phi^L}_{\langle \phi | i \text{ Firing} \rangle} - \underbrace{\int_{\phi^L} \tilde{\mathbf{W}}_i(\phi^L) \phi^L}_{\theta_i} = \langle \phi^L | i \text{ Firing} \rangle - \theta_i. \end{aligned}$$

## Supplementary Note D: Proof of simplest case

The position self-estimate will reach a steady cycle, so we start with a animal at  $\mathbf{r}(t=0) = -\frac{L}{2}$ , having position self-estimate  $\mathcal{R}_0^S$ . The position self-estimate will follow the linearized dynamics, which include terms for path integration as

well as the east and west landmarks.

$$\frac{d\mathcal{R}^S}{dt} = \frac{d\mathbf{r}}{dt} + \omega H_{\text{East}}(\mathbf{r}) (\mathcal{R}_{\text{East}}^L - \mathcal{R}^S) + \omega H_{\text{West}}(\mathbf{r}) (\mathcal{R}_{\text{West}}^L - \mathcal{R}^S)$$

$$H_{\text{East}}(\mathbf{r}) = [\mathbf{r} - (L/2 - L_{\text{Wall}})]_+, \quad H_{\text{West}}(\mathbf{r}) = [(-L/2 + L_{\text{Wall}}) - \mathbf{r}]_+$$

We can assume the position self-estimate will reach a steady cycle such that  $\mathcal{R}^S(t=2\tau) = \mathcal{R}^S(t=0)$ . Defining  $\mathcal{R}_1^S = \mathcal{R}^S(\tau/2)$ ,  $\mathcal{R}_2^S = \mathcal{R}^S(\tau)$ ,  $\mathcal{R}_3^S = \mathcal{R}^S(3\tau/2)$ , we can solve for the the position self-estimate as a piecewise function:

$$\mathcal{R}^S(t) = \begin{cases} \mathcal{R}_1^S + vt & \tau_w < t < \tau - \tau_w \\ (\mathcal{R}_1^S + [L_{\text{Int}}/2]) e^{-\omega(t-\tau_w)} + [\mathcal{R}_{\text{East}}^L + \frac{v}{\omega}] (1 - e^{-\omega t}) & \tau - \tau_w \leq t \leq \tau \\ \mathcal{R}_2^S e^{-\omega(t-\tau_w)} + [\mathcal{R}_{\text{East}}^L - \frac{v}{\omega}] (1 - e^{-\omega t}) & \tau < t \leq \tau + \tau_w \\ \dots \end{cases} \quad (24)$$

Where  $\tau_w = L_{\text{Wall}}/v$ . This yields a set of linear equations:

$$\begin{aligned} \mathcal{R}_1^S &= e^{-\omega\tau_w} \mathcal{R}_0^S + (1 - e^{-\omega\tau_w}) \left[ \mathcal{R}_{\text{West}}^L + \frac{v}{\omega} \right] + [L_{\text{Int}}/2] \\ \mathcal{R}_2^S &= e^{-\omega\tau_w} (\mathcal{R}_1^S + [L_{\text{Int}}/2]) + (1 - e^{-\omega\tau_w}) \left[ \mathcal{R}_{\text{East}}^L + \frac{v}{\omega} \right] \\ \mathcal{R}_3^S &= e^{-\omega\tau_w} \mathcal{R}_2^S + (1 - e^{-\omega\tau_w}) \left[ \mathcal{R}_{\text{East}}^L - \frac{v}{\omega} \right] - [L_{\text{Int}}/2] \\ \mathcal{R}_4^S &= \mathcal{R}_0^S = e^{-\omega\tau_w} (\mathcal{R}_3^S - [L_{\text{Int}}/2]) + (1 - e^{-\omega\tau_w}) \left[ \mathcal{R}_{\text{West}}^L - \frac{v}{\omega} \right] \end{aligned}$$

The average position self-estimate seen by the east landmark comes from two components of piecewise function  $\mathcal{R}^S(t)$ . The first is  $\tau - \tau_w < t < \tau$ , the second is  $\tau < t < \tau + \tau_w$ :

$$\begin{aligned} \bar{\mathcal{R}}_{\text{East}}^S &= \left\langle \mathcal{R}^S(t) \mid \text{Landmark Cell A Firing} \right\rangle = \frac{\int_0^{2\tau} H_{\text{East}}(\mathbf{r}(t)) \mathcal{R}^S(t) dt}{\int_0^{2\tau} H_{\text{East}}(\mathbf{r}(t)) dt} = \frac{1}{2\tau_w} \int_{\tau - \tau_w}^{\tau + \tau_w} \mathcal{R}^S(t) dt = \\ \bar{\mathcal{R}}_{\text{East}}^S &= \frac{1}{2\tau_w} \cdot \left[ \left( \int_0^{\tau_w} \underbrace{(\mathcal{R}_1^S + [L_{\text{Int}}/2]) e^{-\omega t} + [\mathcal{R}_{\text{East}}^L + v] (1 - e^{-\omega t})}_{\tau - \tau_w < t < \tau} \right) + \left( \int_0^{\tau_w} \underbrace{\mathcal{R}_2^S e^{-\omega t} + [\mathcal{R}_{\text{East}}^L - v] (1 - e^{-\omega t})}_{\tau < t < \tau + \tau_w} \right) \right] = \\ \frac{1}{2\tau_w} \cdot &\left[ \left( \mathcal{R}_1^S + [L_{\text{Int}}/2] \right) \left( \frac{1 - e^{-\omega\tau_w}}{\omega} \right) + \left[ \mathcal{R}_{\text{East}}^L + \cancel{v} \right] \left( \tau_w - \frac{1 - e^{-\omega\tau_w}}{\omega} \right) + \mathcal{R}_2^S \left( \frac{1 - e^{-\omega\tau_w}}{\omega} \right) + \left[ \mathcal{R}_{\text{East}}^L - \cancel{v} \right] \left( \tau_w - \left( \frac{1 - e^{-\omega\tau_w}}{\omega} \right) \right) \right] \\ &= \frac{1}{2\tau_w} \cdot \left[ (\mathcal{R}_1^S + [L_{\text{Int}}/2]) \left( \frac{1 - e^{-\omega\tau_w}}{\omega} \right) + \mathcal{R}_{\text{East}}^L \left( \tau_w - \frac{1 - e^{-\omega\tau_w}}{\omega} \right) + \mathcal{R}_2^S \left( \frac{1 - e^{-\omega\tau_w}}{\omega} \right) + \mathcal{R}_{\text{East}}^L \left[ \tau_w - \frac{1 - e^{-\omega\tau_w}}{\omega} \right] \right] \\ &= \mathcal{R}_{\text{East}}^L + \frac{(1 - e^{-\omega\tau_w})}{2\omega\tau_w} \left[ (\mathcal{R}_1^S + [L_{\text{Int}}/2] - \mathcal{R}_{\text{East}}^L) + (\mathcal{R}_2^S - \mathcal{R}_{\text{East}}^L) \right] \end{aligned}$$

Therefore, at equilibrium:

$$\mathcal{R}_{\text{East}}^L = \bar{\mathcal{R}}_{\text{East}}^S = \frac{(\mathcal{R}_1^S + [L_{\text{Int}}/2]) + \mathcal{R}_2^S}{2}$$

There is a translational symmetry to this problem, such that any shifted version of a solution is also a solution. We center

around zero for simplicity, such that  $\mathcal{R}_1^S = -\mathcal{R}_3^S$ ,  $\mathcal{R}_2^S = -\mathcal{R}_4^S$ , and  $\mathcal{R}_{\text{East}}^L = -\mathcal{R}_{\text{West}}^L$ . Combining the above equations and this symmetry gives the steady state solution:

$$\begin{aligned}\mathcal{R}_1^S &= \mathcal{R}_3^S = 0, \\ \mathcal{R}_2^S &= \left( [L_{\text{Int}}/2] + \frac{2[1-e^{-\omega\tau_w}]}{1+e^{-\omega\tau_w}} \left(\frac{v}{\omega}\right) \right) = \left( [L_{\text{Int}}/2] + 2 \tanh(\omega\tau_w/2) \left(\frac{v}{\omega}\right) \right) = -\mathcal{R}_0^S \\ \mathcal{R}_{\text{East}}^L &= \left( [L_{\text{Int}}/2] + \frac{[1-e^{-\omega\tau_w}]}{1+e^{-\omega\tau_w}} \left(\frac{v}{\omega}\right) \right) = \left( [L_{\text{Int}}/2] + \tanh(\omega\tau_w/2) \left(\frac{v}{\omega}\right) \right) = -\mathcal{R}_{\text{West}}^L.\end{aligned}$$

**A. Out of Equilibrium Path-Dependent Shifts and Learning Dynamics .** When the system is out of equilibrium it is convenient to refer to the landmark representations in terms of their deviation from the equilibrium state.

$$\mathcal{R}_{\text{East}}^L = \Delta\mathcal{R}_{\text{East}}^L + (\mathcal{R}_{\text{East}}^L)_{\text{Eq}}, \quad \mathcal{R}_{\text{West}}^L = \Delta\mathcal{R}_{\text{West}}^L + (\mathcal{R}_{\text{West}}^L)_{\text{Eq}}$$

We have the set of linear equations for how much the position self-estimates vary with the landmark position estimates, where we use the shorthand  $\mathcal{R}_i^S = \Delta\mathcal{R}_i^S + (\mathcal{R}_i^S)_{\text{Eq}}$ :

$$\begin{aligned}\Delta\mathcal{R}_1^S &= e^{-\omega\tau_w} \Delta\mathcal{R}_0^S + (1-e^{-\omega\tau_w}) \Delta\mathcal{R}_{\text{West}}^L \\ \Delta\mathcal{R}_2^S &= e^{-\omega\tau_w} \Delta\mathcal{R}_1^S + (1-e^{-\omega\tau_w}) \Delta\mathcal{R}_{\text{East}}^L \\ \Delta\mathcal{R}_3^S &= e^{-\omega\tau_w} \Delta\mathcal{R}_2^S + (1-e^{-\omega\tau_w}) \Delta\mathcal{R}_{\text{East}}^L \\ \Delta\mathcal{R}_4^S &= \Delta\mathcal{R}_0^S = e^{-\omega\tau_w} \Delta\mathcal{R}_3^S + (1-e^{-\omega\tau_w}) \Delta\mathcal{R}_{\text{West}}^L\end{aligned}$$

We can combine the equations for  $\Delta\mathcal{R}_3^S, \Delta\mathcal{R}_2^S$  to get:

$$\begin{aligned}\Delta\mathcal{R}_3^S &= e^{-\omega\tau_w} \Delta\mathcal{R}_2^S + (1-e^{-\omega\tau_w}) \Delta\mathcal{R}_{\text{East}}^L \\ &= e^{-\omega\tau_w} \left[ e^{-\omega\tau_w} \Delta\mathcal{R}_1^S + (1-e^{-\omega\tau_w}) \Delta\mathcal{R}_{\text{East}}^L \right] + (1-e^{-\omega\tau_w}) \Delta\mathcal{R}_{\text{East}}^L \\ &= e^{-2\omega\tau_w} \Delta\mathcal{R}_1^S + \Delta\mathcal{R}_{\text{East}}^L (1-e^{-2\omega\tau_w})\end{aligned}$$

We can express through  $\mathcal{R}_1^S$  in terms of  $\mathcal{R}_3^S$  through symmetry:

$$\begin{aligned}\Delta\mathcal{R}_3^S &= e^{-2\omega\tau_w} \Delta\mathcal{R}_1^S + \Delta\mathcal{R}_{\text{East}}^L (1-e^{-2\omega\tau_w}) \\ \Delta\mathcal{R}_1^S &= e^{-2\omega\tau_w} \Delta\mathcal{R}_3^S + \Delta\mathcal{R}_{\text{West}}^L (1-e^{-2\omega\tau_w})\end{aligned}$$

Plugging one into the other:

$$\begin{aligned}\Delta\mathcal{R}_1^S &= e^{-2\omega\tau_w} \left[ e^{-2\omega\tau_w} \Delta\mathcal{R}_1^S + \Delta\mathcal{R}_{\text{East}}^L (1-e^{-2\omega\tau_w}) \right] + \Delta\mathcal{R}_{\text{West}}^L (1-e^{-2\omega\tau_w}) \Rightarrow \\ \Delta\mathcal{R}_1^S &= (e^{-4\omega\tau_w}) \Delta\mathcal{R}_1^S + e^{-2\omega\tau_w} (1-e^{-2\omega\tau_w}) \Delta\mathcal{R}_{\text{East}}^L + \Delta\mathcal{R}_{\text{West}}^L (1-e^{-2\omega\tau_w}) \Rightarrow \\ (1-e^{-4\omega\tau_w}) \Delta\mathcal{R}_1^S &= e^{-2\omega\tau_w} (1-e^{-2\omega\tau_w}) \Delta\mathcal{R}_{\text{East}}^L + \Delta\mathcal{R}_{\text{West}}^L (1-e^{-2\omega\tau_w}) \Rightarrow \\ (1-e^{-4\omega\tau_w}) \Delta\mathcal{R}_1^S &= (1-e^{-2\omega\tau_w}) \left[ e^{-2\omega\tau_w} \Delta\mathcal{R}_{\text{East}}^L + \Delta\mathcal{R}_{\text{West}}^L \right]\end{aligned}$$

This yields the change in position self-estimate:

$$\Delta\mathcal{R}_1^S = \frac{[e^{-2\omega\tau_w} \Delta\mathcal{R}_{\text{East}}^L + \Delta\mathcal{R}_{\text{West}}^L]}{1+e^{-2\omega\tau_w}} = \Delta\mathcal{R}_{\text{East}}^L + \frac{(\Delta\mathcal{R}_{\text{West}}^L - \Delta\mathcal{R}_{\text{East}}^L)}{1+e^{-2\omega\tau_w}}$$

This allows us to recover the first coefficient related to path-dependent shift. When  $\Delta\mathcal{R}_{\text{East}}^{\text{L}} = -\Delta\mathcal{R}_{\text{West}}^{\text{L}}$ ,

$$\Delta\mathcal{R}_1^{\text{S}} = \Delta\mathcal{R}_{\text{East}}^{\text{L}} + \frac{(\Delta\mathcal{R}_{\text{West}}^{\text{L}} - \Delta\mathcal{R}_{\text{East}}^{\text{L}})}{1+e^{-2\omega\tau_w}} = \Delta\mathcal{R}_{\text{East}}^{\text{L}} \left[ \frac{1+e^{-2\omega\tau_w}-2}{1+e^{-2\omega\tau_w}} \right] = -\Delta\mathcal{R}_{\text{East}}^{\text{L}} \left[ \frac{1-e^{-2\omega\tau_w}}{1+e^{-2\omega\tau_w}} \right] = -\Delta\mathcal{R}_{\text{East}}^{\text{L}} \tanh(\omega\tau_w) \quad (25)$$

We note that when  $\omega\tau_w \rightarrow \infty$ , the path-dependent shift is *exactly* the shift in the estimated position of the landmark last touched  $\mathcal{R}_{\text{West}}^{\text{L}}$ . When  $\omega\tau_w \rightarrow 0$ , the shift goes to 0, as the memory of  $\mathcal{R}_{\text{East}}^{\text{L}}$  is nearly the same as that of  $\mathcal{R}_{\text{West}}^{\text{L}}$ .

**A.1. Learning Timescale Coefficient.** In order to understand the learning dynamics, we must calculate the effect of the *estimated* landmark position on the *estimated* position self-estimate that becomes associated with each landmark:

$$\Delta\bar{\mathcal{R}}_{\text{East}}^{\text{S}} = \mathcal{R}_{\text{East}}^{\text{L}} + \left[ \frac{\Delta\mathcal{R}_1^{\text{S}} + \Delta\mathcal{R}_2^{\text{S}}}{2} - \mathcal{R}_{\text{East}}^{\text{L}} \right] \cdot \frac{(1 - e^{-\omega\tau_w})}{\omega\tau_w}$$

Plugging in:

$$\Delta\mathcal{R}_2^{\text{S}} = \Delta\mathcal{R}_{\text{East}}^{\text{L}} + (\Delta\mathcal{R}_1^{\text{S}} - \mathcal{R}_{\text{East}}^{\text{L}}) e^{-\omega\tau_w}$$

Gives things in terms of  $\mathcal{R}_1^{\text{S}}$ :

$$\Delta\bar{\mathcal{R}}_{\text{East}}^{\text{S}} = \mathcal{R}_{\text{East}}^{\text{L}} + (\mathcal{R}_1^{\text{S}} - \mathcal{R}_{\text{East}}^{\text{L}}) \left[ \frac{1+e^{-\omega\tau_w}}{2} \right] \left[ \frac{(1-e^{-\omega\tau_w})}{\omega\tau_w} \right] = \mathcal{R}_{\text{East}}^{\text{L}} + (\mathcal{R}_1^{\text{S}} - \mathcal{R}_{\text{East}}^{\text{L}}) \left[ \frac{1-e^{-2\omega\tau_w}}{2\omega\tau_w} \right]$$

Plugging in the value of  $\mathcal{R}_1^{\text{S}}$ :

$$\Delta\mathcal{R}_{\text{East}}^{\text{L}} + \frac{(\Delta\mathcal{R}_{\text{West}}^{\text{L}} - \Delta\mathcal{R}_{\text{East}}^{\text{L}})}{1+e^{-2\omega\tau_w}}$$

Gives:

$$\Delta\bar{\mathcal{R}}_{\text{East}}^{\text{S}} = \mathcal{R}_{\text{East}}^{\text{L}} + (\mathcal{R}_{\text{West}}^{\text{L}} - \mathcal{R}_{\text{East}}^{\text{L}}) \left[ \frac{1-e^{-\omega\tau_w}}{2\omega\tau_w(1+e^{-2\omega\tau_w})} \right]$$

Yielding a learning time of:

$$T_{\text{Learning}} = \frac{2\omega\tau_w(1+e^{-2\omega\tau_w})}{1-e^{-\omega\tau_w}} \quad (26)$$

From Eq. 25, we can see that as the landmark cells become stronger, the shifts become stronger, as the animals position self-estimate becomes more heavily weighted toward whichever landmark it most recently saw. From Eq. 26 we see that, as landmark cells become stronger, the learning rate slows down, as landmark cells mostly see their own self-estimates; the contribution to position self-estimate from spatially disjoint landmarks decays quickly after the animal moves into the landmark firing field.

## Supplementary Note E: Periodicity of 2D representation

When 2D attractor dynamics yield a family of steady hexagonal bump patterns, this periodicity can be represented mathematically on the neural sheet as:

$$sss(u_1, u_2) = sss(u_1 + 2\pi, u_2) = sss(u_1 + \pi, u_2 + \sqrt{3}\pi)$$

Where we have defined units on the neural sheet in terms of this periodicity. Therefore, the coordinate  $\vec{\phi}$  specifying a point on the manifold of stable attractor patterns is a periodic variable defined modulo the periodicity of the steady state pattern:

$$\vec{\phi}_A \equiv \vec{\phi}_A + (2\pi, 0) \equiv \vec{\phi}_A + (\pi, \sqrt{3}\pi)$$

## Supplementary Note F: Proof that Learning Dynamics Reduce to a Mechanical Framework

We show how these dynamics can be reduced to a low-dimensional model where the exploration behavior gives an emergent interaction between the learned states of landmark cells, even those with non-overlapping firing patterns. We first need to solve for an animal's position self-estimate as a function of its path history. To do so we first make the bookkeeping substitution:

$$\frac{dR^S}{dt} = \underbrace{\frac{d\mathbf{r}}{dt}}_{\text{Path Integration}} + \underbrace{\omega(\mathbf{r}) [\mathcal{R}^L(\mathbf{r}) - R^S]}_{\text{Landmark Cells}} \quad (27)$$

Where  $\omega(\mathbf{r}) = \sum \omega H_i(\mathbf{r})$  is the combined strength of all landmark cells that fire at  $\mathbf{r}$ , and  $\mathcal{R}^L(\mathbf{r}) = \sum [H_i(\mathbf{r}) R_i^L] / \omega(\mathbf{r})$  is the average position estimate being reinforced at position  $\mathbf{r}$ . As the animal moves around the environment, the position self-estimate will get pushed to the learned positions of landmarks the animal visits, path integrated as the animal moves, and eventually forgotten as the animal orients itself to new landmarks. We can take this basic intuition and turn it into a closed-form equation(Verified in App. A); given any path history  $\mathbf{r}(t)$  the solution for Eq. 27 is:

$$R^S(\mathbf{r}(t), t) = \int_{-\infty}^t \underbrace{[\mathcal{R}^L(\mathbf{r}(t')) + (\mathbf{r}(t) - \mathbf{r}(t'))]}_{\text{Landmark Position Estimate + Path Integration from } t'} \cdot \underbrace{\omega(\mathbf{r}(t')) \cdot \left[ e^{-\int_{t'}^t \omega(\mathbf{r}(t'')) dt''} \right]}_{\text{Memory of time } t'} dt' \quad (28)$$

**Solving for learned position estimates as a function of current landmark position estimates** We now need to compute the mean position-self estimate seen by each landmark cell. We note that for any *individual* path,  $R^S(\mathbf{r}(t), t)$  is linear with respect to  $\mathcal{R}^L(\mathbf{r}')$ . Therefore, averaging over *all* paths  $R^S(\mathbf{r}(t), t)$  that end at  $\mathbf{r}$ , the average  $\bar{R}^S(\mathbf{r})$  is *also* linear with  $\mathcal{R}^L(\mathbf{r}')$ . Therefore, we can construct a matrix equation:

$$\bar{R}^S(\mathbf{r}) = \int S(\mathbf{r}, \mathbf{r}') \left( \omega(\mathbf{r}') [\mathcal{R}^L(\mathbf{r}') + (\mathbf{r} - \mathbf{r}')] \right) d\mathbf{r}',$$

where our matrix entries  $S(\mathbf{r}, \mathbf{r}')$  represent all possible ways the landmark position-estimates at position  $\mathbf{r}'$  contribute to the mean position self-estimate at  $\mathbf{r}$ . As long as the exploration dynamics are reversible, i.e., for any  $\mathbf{r}(t)$ , the reverse path  $\mathbf{r}(-t)$  is equally likely,  $S$  is symmetric ( $S(\mathbf{r}, \mathbf{r}') = S(\mathbf{r}', \mathbf{r})$ ) (Proof in App. B).

To solve for the learning dynamics, we expand  $\omega(\mathbf{r}), \mathcal{R}^L(\mathbf{r})$  to understand the average position self-estimate as a function of position:

$$\bar{R}^S(\mathbf{r}) = \sum_j \int S(\mathbf{r}, \mathbf{r}') H_j(\mathbf{r}') \left( \mathcal{R}_j^L + (\mathbf{r} - \mathbf{r}') \right) d\mathbf{r}'.$$

The mean position self-estimate seen by each landmark cell is then:

$$\bar{R}_i^S = \sum_j \iint_{\mathbf{r}, \mathbf{r}'} H_i(\mathbf{r}) S(\mathbf{r}, \mathbf{r}') H_j(\mathbf{r}') \left( \mathcal{R}_j^L + (\mathbf{r} - \mathbf{r}') \right).$$

Combining this with the angle learning rule gives:

$$\begin{aligned} \frac{d\mathcal{R}_i^L}{dT} &= \bar{\mathcal{R}}_i^S - \mathcal{R}_i^L = \sum_j \iint_{\mathbf{r}, \mathbf{r}'} H_i(\mathbf{r}) S(\mathbf{r}, \mathbf{r}') H_j(\mathbf{r}') [\mathcal{R}_j^L + (\mathbf{r} - \mathbf{r}')] - \mathcal{R}_i^L \\ &= \sum_j \underbrace{\left[ \iint_{\mathbf{r}, \mathbf{r}'} H_i(\mathbf{r}) S(\mathbf{r}, \mathbf{r}') H_j(\mathbf{r}') \right]}_{\mathbf{M}_{ij} \text{ (Definition)}} \mathcal{R}_j^L + \sum_j \underbrace{\left[ \iint_{\mathbf{r}, \mathbf{r}'} H_i(\mathbf{r}) S(\mathbf{r}, \mathbf{r}') H_j(\mathbf{r}') [\mathbf{r} - \mathbf{r}'] \right]}_{\Delta\mathcal{R}_{j \rightarrow i}^S \cdot \mathbf{M}_{ij} \text{ (Definition)}} - \mathcal{R}_i^L \end{aligned}$$

We note that  $\sum_j \mathbf{M}_{ij} = 1$  for all  $i$ <sup>4</sup>; therefore, we can rewrite the above equation as:

$$d\mathcal{R}_i^L/dT = \sum_j \mathbf{M}_{ij} \left( [\mathcal{R}_j^L + \Delta\mathcal{R}_{j \rightarrow i}^S] - \mathcal{R}_i^L \right).$$

Due to the symmetry of  $S$ ,  $\mathbf{M}_{ij} = \mathbf{M}_{ji}$ ,  $\Delta\mathcal{R}_{j \rightarrow i}^S = -\Delta\mathcal{R}_{i \rightarrow j}^S$ . Therefore, the long term dynamics of mapping are equivalent to the first-order dynamics of a set of particles  $i$ , attached by springs of strength  $\mathbf{M}_{ij}$ , having a rest displacement of  $\Delta\mathcal{R}_{j \rightarrow i}^S$ . The spring constant is related to the frequency with which the animal moves between each landmark field, while the rest displacement is related to the distance between the landmark field.

**A. Proof of convolutional integral.** We can check that the solution for the position self-estimate Eq. 28:

$$\mathcal{R}^S(\mathbf{r}(t), t) = \int_{-\infty}^t [\mathcal{R}^L(\mathbf{r}(t')) + (\mathbf{r}(t) - \mathbf{r}(t'))] \cdot \omega(\mathbf{r}(t')) \left[ e^{-\int_{t'}^t \omega(\mathbf{r}(t'')) dt''} \right] dt'$$

Satisfies the dynamics of Eq. 27 :

$$\frac{d\mathcal{R}^S(t)}{dt} = \frac{d\mathbf{r}(t)}{dt} + \omega(\mathbf{r}) [\mathcal{R}^L(\mathbf{r}(t)) - \mathcal{R}^S]$$

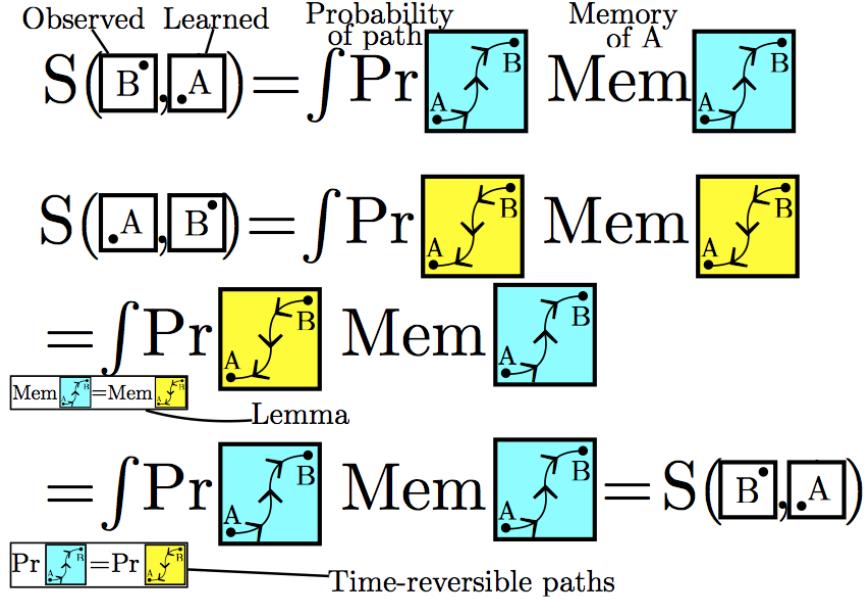
by inspection. We plug Eq. 28 into Eq. 27 to get:

$$\begin{aligned} \frac{d\mathcal{R}^S}{dt} &= \underbrace{[\mathcal{R}^L(\mathbf{r}(t)) + (\mathbf{r}(t) - \mathbf{r}(t))] \left[ \omega(\mathbf{r}(t)) e^{-\int_t^t \omega(\mathbf{r}(t'')) dt''} \right]}_{=\mathcal{R}^L \cdot \omega} + \underbrace{\int_{-\infty}^t \frac{d\mathbf{r}(t)}{dt} \omega(\mathbf{r}(t')) e^{-\int_{t'}^t \omega(\mathbf{r}(t'')) dt''} dt'}_{=\frac{d}{dt}\mathbf{r}} \\ &\quad + \underbrace{\int_{-\infty}^t [\mathcal{R}^L(\mathbf{r}(t')) + (\mathbf{r}(t) - \mathbf{r}(t'))] \cdot \left[ -\omega(\mathbf{r}(t')) \cdot \omega(\mathbf{r}(t)) e^{-\int_{t'}^t \omega(\mathbf{r}(t'')) dt''} \right] dt'}_{=-\omega\mathcal{R}^S} \end{aligned}$$

The underbraced identities are more easily seen by simplifying terms:

$$\begin{aligned} \frac{d}{dt} \mathcal{R}^S &= \underbrace{\left[ \mathcal{R}^L(\mathbf{r}(t)) + (\mathbf{r}(t) - \mathbf{r}(t)) \right]}_{\mathcal{R}^L \cdot \omega} \underbrace{\left[ \omega(\mathbf{r}(t)) e^{-\int_t^t \omega(\mathbf{r}(t'')) dt''} \right]}_0 + \underbrace{\frac{d\mathbf{r}(t)}{dt} \int_{-\infty}^t \omega(\mathbf{r}(t')) e^{-\int_{t'}^t \omega(\mathbf{r}(t'')) dt''} dt'}_{\frac{d\mathbf{r}}{dt}} \\ &\quad - \underbrace{\omega(\mathbf{r}(t)) \cdot \int_{-\infty}^t [\mathcal{R}^L(\mathbf{r}(t')) + (\mathbf{r}(t) - \mathbf{r}(t'))] \cdot \left[ \omega(\mathbf{r}(t')) \cdot e^{-\int_{t'}^t \omega(\mathbf{r}(t'')) dt''} \right] dt'}_{\mathcal{R}^S} \end{aligned}$$

<sup>4</sup>This is due to the fact that shifting *all* position estimates will not change the dynamics.



**Fig. 9.** Sketch of proof in App. B that  $S$  is symmetric for time-symmetric path distributions.

**B. Proof that  $S$  is symmetric for time-symmetric path distributions .** The mean position self-estimate of the animal at position  $\mathbf{r}_B$  is the average self-estimate of all paths that pass  $\mathbf{r}_B$  at time  $t = 0$ . (We pick  $t = 0$  for mathematical convenience).

$$\bar{\mathcal{R}}^S(\mathbf{r}_B) = \int \mathcal{D}\mathbf{r}(t) \Pr(\mathbf{r}(t)) \delta(\mathbf{r}(0) - \mathbf{r}_B) \mathcal{R}^S(\mathbf{r}(0), t = 0),$$

To avoid clutter, use the shorthand:

$$F(\mathbf{r}, t, t') = [\theta(\mathbf{r}(t')) + (\mathbf{r}(t) - \mathbf{r}(t'))] \cdot \omega(\mathbf{r}(t')), \quad \text{Mem}(\mathbf{r}, t, t') = e^{-\int_{t'}^t \omega(\mathbf{r}(t'')) dt''},$$

And decompose this into contributions from different past times  $t'$ .

$$\bar{\mathcal{R}}^S(\mathbf{r}_B) = \int \mathcal{D}\mathbf{r}(t) \Pr(\mathbf{r}(t)) \delta(\mathbf{r}(0) - \mathbf{r}_B) \left[ \int_{-\infty}^0 F(\mathbf{r}, t = 0, t') \text{Mem}(\mathbf{r}, t = 0, t') dt' \right]$$

Reshuffling the order of integration and breaking things down further into contributions of  $\mathbf{r}_A = \mathbf{r}(t')$

$$\bar{\mathcal{R}}^S(\mathbf{r}_B) = \int_{-\infty}^0 dt' \int d\mathbf{r}_A \int \mathcal{D}\mathbf{r}(t) \Pr(\mathbf{r}(t)) \delta(\mathbf{r}(0) - \mathbf{r}_B) \delta(\mathbf{r}(t') - \mathbf{r}_A) [F(\mathbf{r}, t = 0, t') \text{Mem}(\mathbf{r}, t = 0, t')]$$

Because we have assumed the statistics of the animal trajectories  $\mathbf{r}(t)$  will be time-reversal symmetric, the reverse, time shifted path  $\tilde{\mathbf{r}}_{\text{rev}}(t) = \mathbf{r}(t' - t)$  is equally likely. We therefore apply the symmetrization procedure:

$$2 \cdot \bar{\mathcal{R}}^S(\mathbf{r}_B) = \int_{-\infty}^0 dt' \int d\mathbf{r}_A \int \mathcal{D}\mathbf{r}(t) \Pr(\mathbf{r}(t)) \dots \quad (29)$$

$$[\delta(\mathbf{r}(0) - \mathbf{r}_B) \delta(\mathbf{r}(t') - \mathbf{r}_A) F(\mathbf{r}, 0, t') \text{Mem}(\mathbf{r}, 0, t')] + [\delta(\tilde{\mathbf{r}}_{\text{rev}}(0) - \mathbf{r}_B) \delta(\tilde{\mathbf{r}}_{\text{rev}}(t') - \mathbf{r}_A) F(\tilde{\mathbf{r}}_{\text{rev}}, 0, t') \text{Mem}(\tilde{\mathbf{r}}_{\text{rev}}, 0, t')] \quad (30)$$

We note that  $\text{Mem}(\tilde{\mathbf{r}}_{\text{rev}}, 0, t') = \text{Mem}(\mathbf{r}, 0, t')$ , and that:

$$\delta(\tilde{\mathbf{r}}_{\text{rev}}(0) - \mathbf{r}_B)\delta(\tilde{\mathbf{r}}_{\text{rev}}(t') - \mathbf{r}_A)F(\tilde{\mathbf{r}}_{\text{rev}}, 0, t') = \delta(\mathbf{r}(t') - \mathbf{r}_B)\delta(\mathbf{r}(0) - \mathbf{r}_A)F(\mathbf{r}, 0, t')$$

Therefore, we can simplify Eq. 30 to:

$$2 \cdot \bar{\mathcal{R}}^S(\mathbf{r}_B) = \int d\mathbf{r}_A \int_{-\infty}^0 dt' \int D\mathbf{r}(t) \Pr(\mathbf{r}(t)) \cdot \quad (31)$$

$$\text{Mem}(\mathbf{r}, 0, t')[\theta(\mathbf{r}_A) + \mathbf{k} \cdot (\mathbf{r}_B - \mathbf{r}_A)] \cdot \omega(\mathbf{r}_A) [\delta(\mathbf{r}(0) - \mathbf{r}_B)\delta(\mathbf{r}(t') - \mathbf{r}_A) + \delta(\mathbf{r}(0) - \mathbf{r}_A)\delta(\mathbf{r}(t') - \mathbf{r}_B)] = \quad (32)$$

$$2 \int d\mathbf{r}_A S(\mathbf{r}_A, \mathbf{r}_B) [\theta(\mathbf{r}_A) + \mathbf{k} \cdot (\mathbf{r}_B - \mathbf{r}_A)] \cdot \omega(\mathbf{r}_A) \quad (33)$$

Where our matrix entries:

$$2S(\mathbf{r}_A, \mathbf{r}_B) = \int_{-\infty}^0 dt' \int D\mathbf{r}(t) \Pr(\mathbf{r}(t)) \text{Mem}(\mathbf{r}, 0, t') [\delta(\mathbf{r}(0) - \mathbf{r}_B)\delta(\mathbf{r}(t') - \mathbf{r}_A) + \delta(\mathbf{r}(0) - \mathbf{r}_A)\delta(\mathbf{r}(t') - \mathbf{r}_B)]$$

are symmetric with respect to the swapping of  $\mathbf{r}_B, \mathbf{r}_A$ .

This proof assumes uniform density of animal positions with uniform areas and strengths of each landmark cell. The proof can be generalized beyond these constraints by making effective particles corresponding to certain landmarks more “massive”, but here we present the simpler proof in the interest of clarity.

## Supplementary Note G: Simulations

**A. Exploration.** In our simulations, we discretize space onto a grid. For simplicity, we have the animal follow diffusive dynamics, implemented through a random walk; at every time step, the animal moves to one of four neighboring cells; any move which would take the animal outside the box is prohibited. The animal has a position self-estimate  $\mathcal{R}^S(t)$  as well as an attractor state  $\vec{\phi}_A(t)$ , which undergoes discrete path-integration at every time step:

$$\begin{aligned} \mathcal{R}^S(t + \Delta t) &\rightarrow \mathcal{R}^S(t) + \Delta\mathbf{r}(t), \\ \vec{\phi}_A(t + \Delta t) &\rightarrow \vec{\phi}_A(t) + \overset{\leftrightarrow}{K}\Delta\mathbf{r}(t) \end{aligned}$$

Afterwards, the position self-estimate is pulled towards the position estimates of any landmark cells which are firing:

$$\begin{aligned} \mathcal{R}^S(t + \Delta t) &\rightarrow \mathcal{R}^S(t + \Delta t) + \left( \omega(\mathbf{r}) \Delta t [\mathcal{R}^L(\mathbf{r}) - \mathcal{R}^S(t + \Delta t)] \right) \cdot \Delta t \\ \vec{\phi}_A(t + \Delta t) &\rightarrow \vec{\phi}_A(t + \Delta t) + \left( \sum_i H_i(\mathbf{r}(t)) \sum_{\vec{\phi}^L} \tilde{W}_i(\vec{\phi}^L) \text{Force}_A(\vec{\phi}_A - \vec{\phi}^L) \right) \cdot \Delta t, \end{aligned}$$

Where  $\vec{\phi}^L$  is discretized into a  $15 \times 15$  grid so that  $\tilde{W}_i(\vec{\phi}^L)$  can be represented as an array.

We set the timescale of animal motion to be 1

$$\Delta t = |\Delta\mathbf{r}|^2 \cdot D$$

Which removes dependence on the discretization size.

**B. Learning.** The learned states are initialized to their firing field center of masses. At every learning epoch T, the simulated animal is placed in the box with an initial position and position self-estimate and explores to get good

statistics.  $\bar{\mathcal{R}}^S(\mathbf{r})$ , is logged, and at the end of each learning epoch, the position estimate of each landmark cell  $i$  is updated to be the average position self-estimate when the landmark cell is firing.

$$\begin{aligned}\mathcal{R}_{i,T+1}^L &\rightarrow \bar{\mathcal{R}}_{i,T}^S, \\ \hat{\mathbf{W}}_i\left(\vec{\phi}^L\right) &\rightarrow \Pr(\vec{\phi}_A(t) = \vec{\phi}^L | i \text{ Firing})\end{aligned}$$

Each of these will converge after a handful of learning epochs; in practice, we use twenty.

**C. Simulation of Square and Bent Environments.** Landmark cell firing fields are heterogeneous; while some are distributed across an entire border; two replicate this distribution we have two types of landmark cells in our model. 1) Landmark cells having uniform wall-length firing field, with a width of 10cm, for example  $H(x,y) = e^{-\left(\frac{x-x_{\text{wall}}}{5\text{cm}}\right)^2}$  for a landmark cell on the west wall. 2) More localized, overlapping, firing fields along each wall. Each firing field is a 5 cm  $\times$  10 cm half-ellipse of along a particular wall; i.e.  $H(x,y) = e^{-\left(\frac{y-y_0}{10\text{cm}}\right)^2 - \left(\frac{x-x_{\text{wall}}}{5\text{cm}}\right)^2}$  for a landmark field along the EW wall with center  $y_0$ . Each type of landmark cell is evenly distributed along each wall, with the total strength and number set such that total firing strength of localized and non-localized cells is the same, and their combined strength leads to a forgetting time of  $\omega = 8\text{Hz}$  along each wall.

Grid spacing is chosen to be 30 cm for square environments(1  $\times$  1 meter); We set the diffusive constant D to be (10 cm) $^2$ /s such that it takes an animal  $\sim$ 100 seconds to traverse the width of the environment.

Grid spacing is 50 cm with a 7° offset for the first trapezoidal environment (1.9  $\times$  .8 meters, same geometry as (11)); We set the diffusive constant D to be (20 cm) $^2$ /s such that it takes an animal  $\sim$ 100 seconds to traverse the length of the environment.

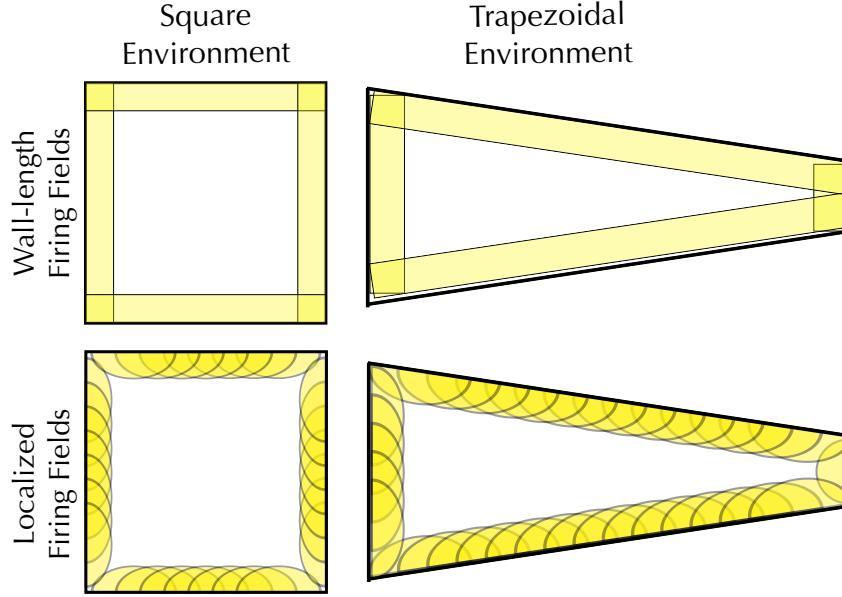
Grid spacing is 50 cm with a 7° offset for the second trapezoidal environment (1.9 meters long. Two straight walls with lengths of .12 meters, .6 meters, with diagonal walls starting 1 meter from the smaller straight wall (14° angle); We set the diffusive constant D to be (20 cm) $^2$ /s such that it takes an animal  $\sim$ 100 seconds to traverse the length of the environment.

The angular offset breaks the symmetry of the trapezoidal environments, yielding bending, but is not required to yield path-dependent shifts.

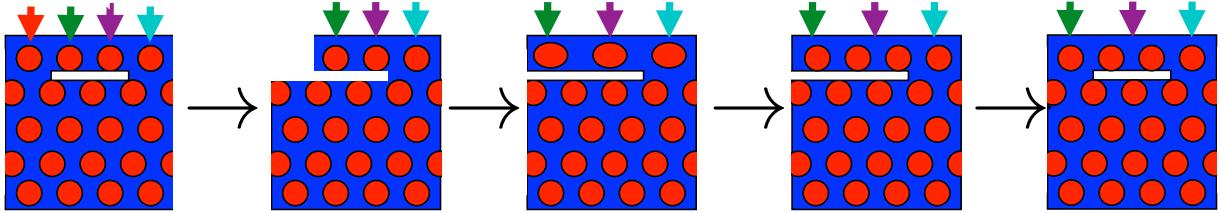
**D. Simulation of Topological Environments.** In order to support topological defects, the environment must be filled with rich, localized landmark cues. In order for the environment to support topological defects, cues must be rich and localized, leading to uniformly distributed landmark firing fields. To model this, we have uniformly localized landmark fields, with  $H(x,y) = e^{-\left(\frac{y-y_0}{10\text{cm}}\right)^2 - \left(\frac{x-x_0}{10\text{cm}}\right)^2}$ , arranged at a density such that their combined strength leads to a forgetting time of 1Hz throughout the environment. The environment was 1.8 meters  $\times$  1 meter, with a center rectangular section of 1.3  $\times$  .8 meters removed. K is chosen to yield a grid spacing of 60cm.

Essential ingredients to achieve topological defects are:

- A “donut-shaped” environment, which can support the topological defect.
- An environment rich in localized, strong landmark cues.
- The larger the environment is, the less deformation it has to support per unit distance, i.e. if an environment is 3 firing fields wide, a topological defect must modify the grid spacing by 33%; if the environment was 5 firing fields wide, the grid spacing would only need to be modified by 20%.
- During the “winding” procedure, the animal cannot acclimate to the intermediate environment for too long; if it fully learns the intermediate environment, the winding procedure will not work (Fig. 11).



**Fig. 10.** Schematic of the distribution of landmark cells for simulations of square and trapezoidal environments. To model a heterogeneous distribution of landmark cell degrees of localization, we include both landmark cells which fire uniformly along a boundary, as well as semi-elliptical landmark cells which are localized to a section of a boundary.



**Fig. 11.** Example of topological defect failing to form due to learning. If the winding procedure is done too slowly, the animal will learn the deformed geometry(Third box → Fourth Box), removing the topological effect.

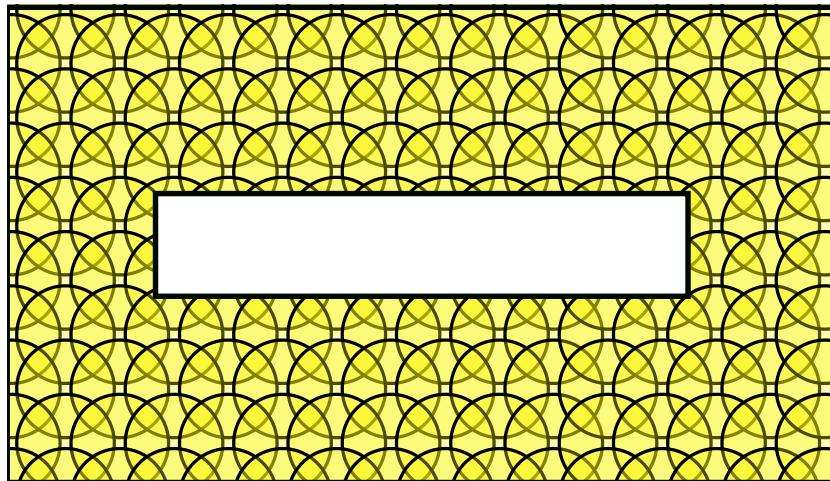
**E. Force law and visualization.** The Grid Cell pattern is visualized by using a truncated parabolic firing rate  $\max\left[1 - \left(\frac{\vec{\phi}_A - \vec{\phi}_0}{D}\right)^2, 0\right]$ , where the field width D is chosen to be  $2\pi/5$ .

The force law chosen is a truncated sin function:

$$\text{Force}_A(\vec{\phi}^L - \vec{\phi}_A) = \begin{cases} -\sin\left(\left|\vec{\phi}^L - \vec{\phi}_A\right|\right) \cdot \frac{\vec{\phi}^L - \vec{\phi}_A}{\left|\vec{\phi}^L - \vec{\phi}_A\right|} & \left|\left(\vec{\phi}^L - \vec{\phi}_A\right)\right| < \pi \\ 0 & \left|\left(\vec{\phi}^L - \vec{\phi}_A\right)\right| \geq \pi \end{cases} \quad (34)$$

We choose this function because it has the correct qualitative features. In addition, in experimental data, the width of a firing rate peak is on the order of the spacing between two firing peaks; this prohibits a force law which is much more short-ranged than this (App. A).

## Topological Environment (Rich Localized Cues)



**Fig. 12.** Schematic of the distribution of landmark cells for simulations the topological environment; cues are densely and uniformly localized throughout the arena.

## Supplementary Note H: Experimental Methods

Data included a subset of published neural recordings previously presented in Hardcastle et al., 2017, Hardcastle et al., 2015. Briefly, mice explored a square box while foraging for chocolate cheerios sprinkled on the floor. During each recording, neural signals from medial entorhinal cortex were recorded and subsequently clustered into distinct neurons. A grid score was computed for each cell following Langston et al., 2010. Cells above a threshold of .4 were considered grid cells. Each grid cell in the dataset was recorded after an average of 28 (data selected from Hardcastle et al., 2015) or 20 (data selected from Hardcastle et al., 2017) exposures to the recording environment.

**A. Estimation of path-dependent shifts.** We examined how grid firing patterns change depending on which landmark (one of four borders in an open 1 meter box) an animal most recently encountered. To control for the effect of head direction and running speed, we preprocessed the data by translating

$$\mathbf{r}(t) \rightarrow \mathbf{r}(t) + 1\text{cm} \times \hat{\mathbf{H}}(t),$$

where  $\hat{\mathbf{H}}(t)$  is a unit vector representing the animal's head direction. This is to avoid artifacts related to tracking; a purely position-dependent firing rate model depends on *some* part of the animal's body, which unlikely to be exactly the position of the tracking diode.

**A.1. Path conditioned rate maps:** . We constructed maps of firing rate as a function of spatial position conditioned on the animal having last touched the North wall more recently than it touched the South Wall, etc. An animal was defined to have “touched” a wall when the head-tracking diodes came within 10cm of the wall. Varying this distance did not significantly effect our results. We avoid any sort of smoothing to prevent artifacts which might show up an experimental signature; as such, the bin size of the computed  $s_{\text{GC}}^{\mathcal{C}}(\mathbf{r})$  is  $5\text{cm} \times 5\text{cm}$ , and each individual trial leaves many bins for which  $s_{\text{GC}}^{\mathcal{C}}(\mathbf{r})$  is not defined. We can create finer-grained cross-correlograms by choosing bin sizes of  $\frac{5}{3}\text{ cm}$ , and smoothing in the manner of (7), but these maps are not used for showing statistical significance. A sort of cross-correlation was taken, using the “angle” between two path-conditioned rate maps.

$$\mathbf{C}_{\text{GC}}^{C_1 C_2}(\Delta \mathbf{r}) = \frac{\| s_{\text{GC}}^{C_1}(\mathbf{r} + \Delta \mathbf{r}) s_{\text{GC}}^{C_2}(\mathbf{r}) \|}{\| s_{\text{GC}}^{C_1}(\mathbf{r} + \Delta \mathbf{r}) \| \| s_{\text{GC}}^{C_2}(\mathbf{r}) \|}$$

Where the mean firing rate is subtracted, and the inner product is only calculated using bins where there is data.

To show significance, we calculate

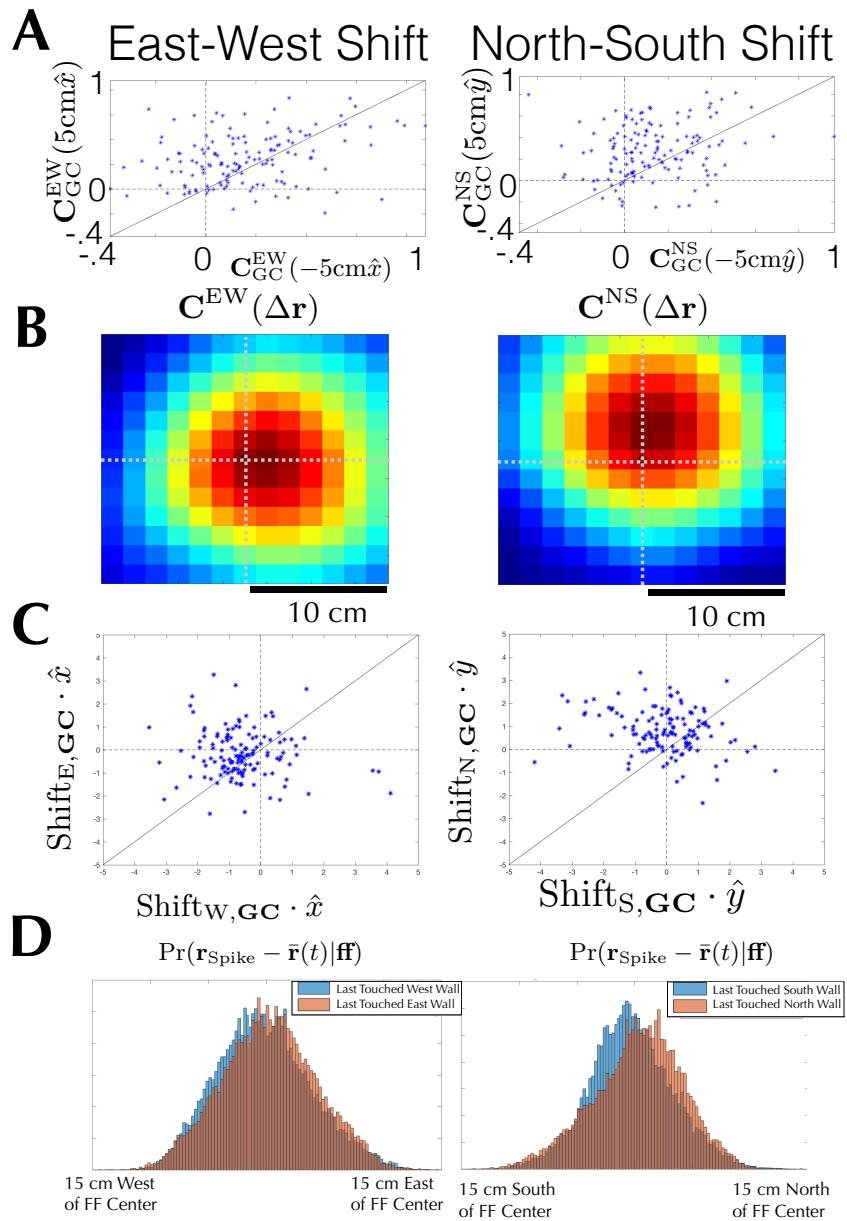
$$\mathbf{C}_{\text{GC}}^{\text{EW}}(5\text{cm}\hat{x}) - \mathbf{C}_{\text{GC}}^{\text{EW}}(-5\text{cm}\hat{x}), \quad \mathbf{C}_{\text{GC}}^{\text{NS}}(5\text{cm}\hat{y}) - \mathbf{C}_{\text{GC}}^{\text{NS}}(-5\text{cm}\hat{y})$$

And show that the patterns are shifted towards whichever wall the animal last touched for both the EW Walls ( $\mathbf{C}_{\text{GC}}^{\text{EW}}(5\text{cm}\hat{x}) - \mathbf{C}_{\text{GC}}^{\text{EW}}(-5\text{cm}\hat{x}) > 0$ ,  $P = 1.5 \cdot 10^{-5}$ , Binomial Test,  $P = 1.5 \cdot 10^{-5}$ , Sign-Flip Test), and the NS walls ( $\mathbf{C}_{\text{GC}}^{\text{NS}}(5\text{cm}\hat{y}) - \mathbf{C}_{\text{GC}}^{\text{NS}}(-5\text{cm}\hat{y})$ ,  $P = 10^{-7}$ , Binomial Test,  $P = 10^{-7}$ , Sign-Flip Test).

**A.2. Spike Displacement:** . Starting from an adaptively smoothed firing rate map, we calculate firing field centers. For each firing field center, we gather positions of spikes recorded in that neighborhood, comparing the average spike position in that neighborhood with the firing center.

$$\begin{aligned} \text{Shift}_{\mathcal{C}, \text{GC}, \mathbf{ff}} &= \langle \mathbf{r}_{\text{Spike}} - \mathbf{r}_{\mathbf{ff}} | \mathcal{C}, \mathbf{r}_{\text{Spike}} \in \mathbf{ff} \rangle \\ &\quad - \langle \mathbf{r}(t) - \mathbf{r}_{\mathbf{ff}} | \mathcal{C}, \mathbf{r}(t) \in \mathbf{ff} \rangle \end{aligned}$$

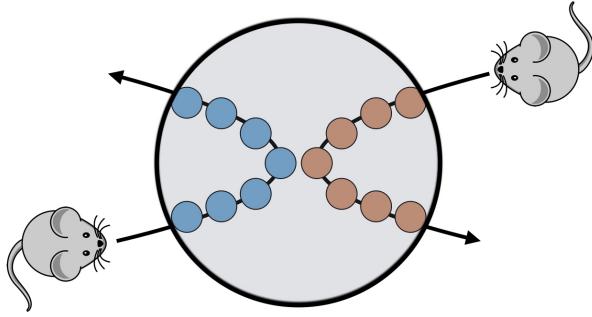
For each path condition  $\mathcal{C}$ , and each firing field center  $\mathbf{ff}$ , we calculate the average spike position  $\mathbf{r}_{\text{Spike}}$  within that firing field, and subtract the average *mouse* position  $\mathbf{r}(t)$  within that firing field. The animal's position within the firing field is subtracted to eliminate any systematic biases that might come from the animal trajectory rather than the actual neural activity (Fig. 14).



**Fig. 13.** **A)** Path-dependent shifts demonstrated by Cross-Correlograms of individual grid cells. Most cells fall on the upper left of the plots, showing that the patterns tend to be shifted towards whichever wall the animal last touched for both the EW Walls ( $P = 1.5 \cdot 10^{-5}$ , Binomial Test,  $P = 1.5 \cdot 10^{-5}$ , Sign-Flip Test), and the NS walls ( $P = 10^{-7}$ , Binomial Test,  $P = 10^{-7}$ , Sign-Flip Test). **B)** The path-dependent shifts is best visualized through the Cross-Correlogram averaged over all grid cells. **C)** Path-dependent shifts demonstrated by Cross-Correlograms of individual grid cells. Most cells fall on the upper left of the plots, showing that the patterns are shifted towards whichever wall the animal last touched for both the EW Walls ( $P = 3 \cdot 10^{-4}$  Binomial Test,  $P = 2 \cdot 10^{-2}$  Sign-Flip Test), and the NS walls ( $P = 10^{-5}$  Binomial Test,  $P = 10^{-5}$  Sign-Flip Test). **D)** The path-dependent shifts is best visualized through a histogram of individual spike displacements.

We calculate the path-dependent shift of an individual grid cell as the average shift of all firing fields in the center:

$$\text{Shift}_{C,\text{GC}} = \sum_{\text{ff}} \text{Shift}_{C,\text{GC},\text{ff}}$$



**Fig. 14.** Schematic of the motivation for subtracting mouse position. An animal is most likely to be closest to the last wall it touched; if the mean animal position was *not* subtracted from the mean spike position, this would yield a path-dependent shift in spike positions purely dependent on animal trajectory rather than neural activity.

To show significance, for each cell, we calculate

$$(\text{Shift}_{E,\text{GC}} - \text{Shift}_{W,\text{GC}}) \cdot \hat{x}, (\text{Shift}_{N,\text{GC}} - \text{Shift}_{S,\text{GC}}) \cdot \hat{y}$$

showing that the patterns are shifted towards whichever wall the animal last touched for both the EW Walls ( $[\text{Shift}_{E,\text{GC}} - \text{Shift}_{W,\text{GC}}] \cdot \hat{x} > 0$ ,  $P = 3 \cdot 10^{-4}$  Binomial Test,  $P = 2 \cdot 10^{-2}$  Sign-Flip Test), and the NS walls ( $[\text{Shift}_{N,\text{GC}} - \text{Shift}_{S,\text{GC}}] \cdot \hat{y} > 0$ ,  $P = 10^{-5}$  Binomial Test,  $P = 10^{-5}$  Sign-Flip Test). We perform both binomial tests, which only depend on the sign of  $(\text{Shift}_{E,\text{GC}} - \text{Shift}_{W,\text{GC}})$ , and magnitude-weighted sign-flip tests, for completeness.