# How to build a cognitive map: insights from models of the hippocampal formation

**James C.R. Whittington**[1,2,†,*], **David McCaffary**[2,†], **Jacob J.W. Bakermans**[2], and **Timothy E.J. Behrens**[2,3,4]

[1]Department of Applied Physics, Stanford University, Stanford, CA, USA
[2]Wellcome Centre for Integrative Neuroimaging, University of Oxford, Oxford, OX3 9DU, UK
[3]Wellcome Centre for Human Neuroimaging, University College London, London, WC1N 3AR, UK
[4]Sainsbury Wellcome Centre for Neural Circuits and Behaviour, University College London, London W1T 4JG, UK
[*]corresponding author(s): `jcrwhittington@gmail.com`
[†]these authors contributed equally to this work

## ABSTRACT

Learning and interpreting the structure of the environment is an innate feature of biological systems, and is integral to guiding flexible behaviours for evolutionary viability. The concept of a *cognitive map* has emerged as one of the leading metaphors for these capacities, and unravelling the learning and neural representation of such a map has become a central focus of neuroscience. While experimentalists are providing a detailed picture of the neural substrate of cognitive maps in hippocampus and beyond, theorists have been busy building models to bridge the divide between neurons, computation, and behaviour. These models can account for a variety of known representations and neural phenomena, but often provide a differing understanding of not only the underlying principles of cognitive maps, but also the respective roles of hippocampus and cortex. In this Perspective, we bring many of these models into a common language, distil their underlying principles of constructing cognitive maps, provide novel (re)interpretations for neural phenomena, suggest how the principles can be extended to account for prefrontal cortex representations and, finally, speculate on the role of cognitive maps in higher cognitive capacities.

## Introduction

Since the 1950s, the hippocampal formation has been implicated in numerous different functions, ranging from episodic memory to spatial and abstract cognition[1–5]. In this time, neuroscientists have attempted to characterise, and provide normative explanations for, the neural representations supporting such functions. Nowhere has this approach proved more fruitful than in the spatial domain, where a variety of cell types, including hippocampal place cells and entorhinal grid cells, provide an appealing neural instantiation of Tolman's (and Turner's) cognitive map[3,4,6,7] (Figure 1a).

Cognitive maps were originally proposed as internal neural representations affording flexible behaviour, such as planning routes or taking never-before-seen shortcuts[6–8]. More recent descriptions formalised this view with the key concept of **generalisation**[2,3,9]. Here, the fundamental role of cognitive maps is to organise knowledge, facilitating generalisation of this knowledge to novel experiences, and thus enabling the rapid inference from sparse observations which characterises biological intelligence[10,11]. Psychologists have thought similarly, both with schemas[12] (a mental framework for understanding new information), and with 'learning to learn'[13] (learning underlying rules of tasks that permit more efficient learning for each new task instantiation). While all these concepts are broad, encompassing domains from social to logical cognition[6], most neural evidence for a cognitive map is grounded in studies of space[3,14].

Recent experimental results, however, increasingly suggest deep parallels between spatial cognition and abstract, non-spatial reasoning[9] (Figure 1b). For instance, hippocampal place cells, which fire with remarkable precision when the animal is in one location in space, also code for one 'place' in sound frequencies[15] when sound frequency is an important component of the given task, or one 'place' in abstract spaces mapped via value[16] or integrated sensory evidence[17]. Similarly, the characteristic hexagonal firing pattern of entorhinal grid cells, discovered in the context of physical space[4], is also found when animals navigate abstract spaces[18–21]. For example, in fMRI, human entorhinal cortex (and medial prefrontal cortex; mPFC)[18,20,21] and monkey mPFC[19] display a hexagonally symmetric pattern when stimuli varying along two abstract dimensions (such as neck and leg length of birds[18], odours[20], social hierarchies[21], or reward probability and value[19]) are presented. These parallels in representation suggest the mechanism for constructing the spatial cognitive map might, in fact, be an instance of a more general coding mechanism capable of building abstract cognitive maps covering any domain.

This presents the exciting and novel opportunity to understand how the brain represents these apparently divergent domains of cognition in the same way. Developing such an understanding, however, requires a formalism connecting physical and abstract space[9]. In recent years, many models of the hippocampal formation have attempted to do this, providing explanations of neural data and offering falsifiable predictions. While greatly informative, these models differ in their focus and the language of their formalism, obscuring the overall direction and vast potential of this work. The aim of this Perspective is to clarify the common theory underlying these models [i], while providing novel results offering normative explanations for a range of old and new neural phenomena. We conclude with a prospective account, speculating how far these models might take us into understanding the neural representations of higher-order cognitive domains, such as language, logical operators, and mathematics, thereby providing a pathway towards cognitive maps as Tolman envisaged - the basis of reasoning across all domains of cognition.

## The cognitive mapping problem

Cognitive maps organise knowledge to afford flexible behaviour[3, 6, 9, 30]. Affording *behaviour* means that the cognitive map must contain information relevant to downstream behavioural tasks. Affording *flexibility* means the map must 1) afford new behaviours in the face of new challenges, and 2) be built as fast as possible, ideally immediately, for any new world - a concept known as *systematic generalisation* in the machine learning literature[31]. The aim for cognitive maps, then, is to learn as much as possible ahead of time, so online learning and computations are minimised - in essence, *learning so that you do not have to learn*. In order to achieve this, there are some requirements and desiderata for the neural representations of the cognitive map. These computational considerations have led to models of the system which have had many recent successes in predicting neuronal representations. Here we describe these computational considerations and explain (in linked boxes) the models relevant to each. We aim to provide a clear conceptual understanding of the interlinked ideas. .

### Reinforcement learning and planning

To afford successful behaviour, cognitive maps must represent *state* (a particular configuration of the world). Knowing when to turn right or left while driving requires an understanding of the orientation of the steering wheel, how far pedals are pushed, the curvature of roads, content of road signs, the location of other vehicles, etc. - turning left when the road bends right can be problematic, after all. Reinforcement learning (RL[32]) is a formalism of this concept: actions are taken based on the current world state (for instance, turning right when the road bends right). Representing the entire world state (in this analogy, not only the steering wheel but also the positions of planets, etc.) is often infeasible, as it can contain information along countless dimensions - and often dimensions irrelevant for solving the current task. Not only is this problematic for representational capacity, but it also impedes the efficiency of learning (it should not be necessary to re-learn how to drive when moving from a red to a black car). This effect, known as the 'curse of dimensionality'[33], can be mitigated by an appropriate state *abstraction* (for instance, ignoring colour in the case of cars). Learning, or attending to, the appropriate abstraction is a central issue of the cognitive mapping problem[22, 34, 35].

Classic (or *model free*) RL learns the value of states, or which actions are good in which states, and therefore requires no knowledge of how states relate to each other. While this is provably optimal in the long term[36], value-based learning is often inflexible and slow to learn[32]. Knowing relationships between states - that is, knowing the state-space structure - however, lets you play different games. Now you can flexibly *plan* routes between any start and any goal state. For instance, taking never-before-experiences routes home[37] (that is, even if those states have no values attached to them) should one's normal route home be blocked. Unfortunately, traditional planning mechanisms (such as tree-search) are computationally costly, but alternatives do exist[38–40]. More broadly, with a clever representation of the state-space (see next section), the cost of planning can be reduced, or even completely avoided. This is a powerful way to formalise the central goal of cognitive maps: building maps that help solve problems in representation, not by exhaustive computation.

### Space as a state-space

To understand what this means, let us consider the state-space of physical space. Here, rather than representing sensory configurations, the state-space comprises physical locations like a literal map. This abstraction alone clearly profoundly helps the spatial planning problem. However, location can be represented in a variety of ways; for example, an unique identifier for each physical location, or (x,y)-coordinates. The choice of which representation to use has major consequences. For example, to find the shortest route: In the former, you need search through a series of neighbours. In the latter, you can compute a vector by subtracting start from end representations. Consider adding a new location. In the former, a new identifier (cell!) is required along with new relationships (synapses!) to neighbouring identifiers. In the latter, nothing new is required - since (x,y) naturally

---

[i]We note that there are many other theoretical accounts of cognitive maps which do not address this issue of representation learning[22–24]. While these models have provided mechanistic insights, they are not discussed in detail here.

extends to new locations. These two representation types are analogous to place and grid cells in the hippocampal formation. Individual place cells code for unique locations, and thus *new* place cells are required for new locations, while grid cells afford
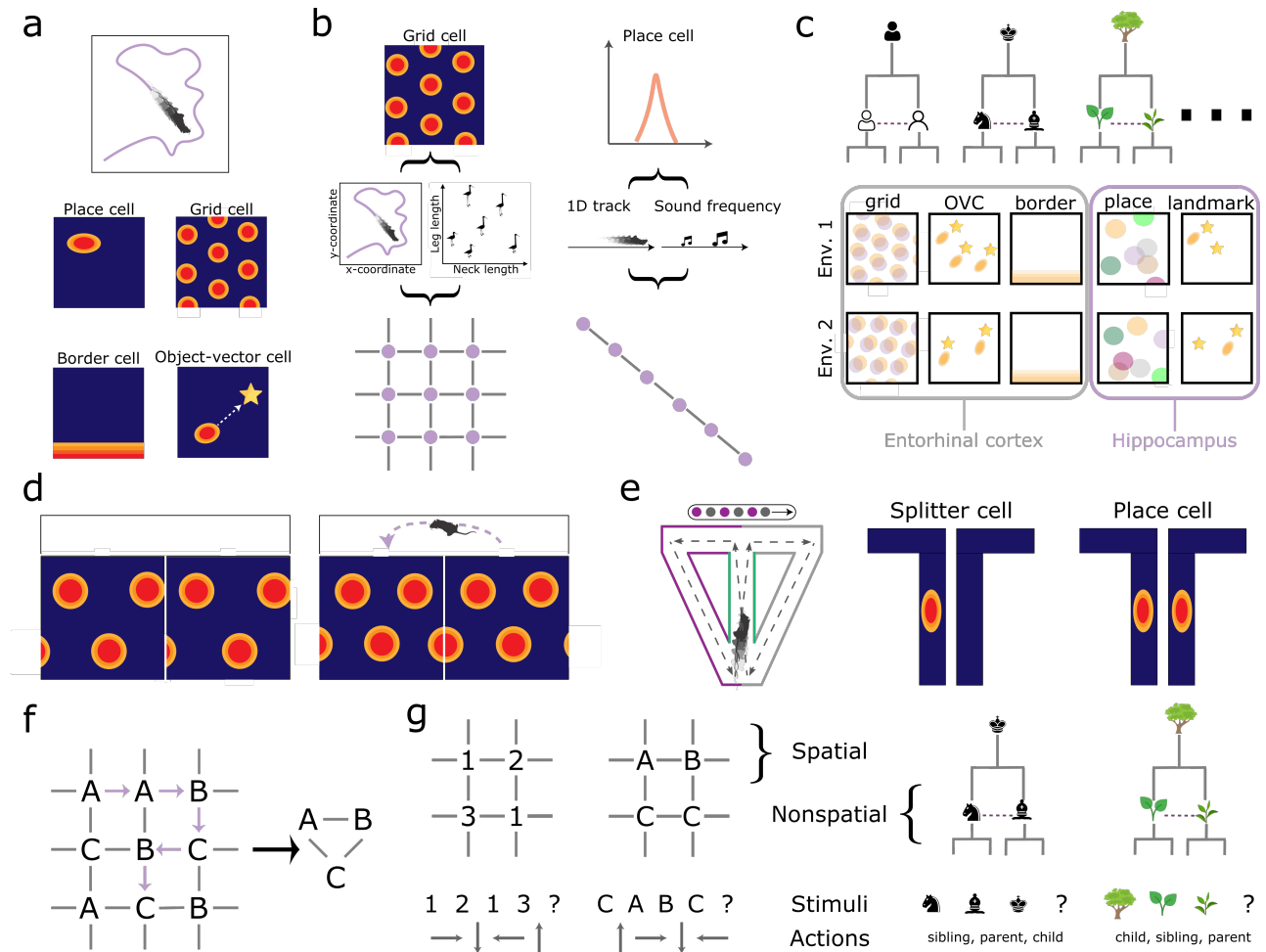


**Figure 1. The cognitive mapping problem: generalisation and latent states. (a)** When navigating naturalistic environments, a range of cell representations are found in the cognitive map of the hippocampal-entorhinal system. Many of these representations were discovered in the context of the spatial cognitive map in rodents[4, 25]. **(b)** Recent evidence, however, has implicated these same representations in the coding of abstract or conceptual spaces (e.g. 'bird space'[18] or sound frequencies[15]), with subsequent theoretical accounts suggesting a single coding mechanism underlies both physical and conceptual spaces[9, 26, 27]: understanding the relational structure (e.g. via *graphs*) of the world. **(c)** Understanding relational structure abstractly affords generalisation, since sensory particularities can differ across instances of the same underlying structure[9]. Cell representations of the entorhinal cortex map relational spaces and *generalise* across environments more so than hippocampal cells. **(d)** Since the sensory world is aliased, representations must be *latent* (not a simple function of the current observation). Rodents exhibit latent state representations when they traverse two sensorially identical rooms[28]. Initially, an identical grid cell code represents both rooms, though when the animal realises these two rooms are connected by a corridor, a *global* grid cell code predominates; the latent representations separate out states with differing sensory futures. **(e)** Latent states are not just in space. In a T-maze alternation task[29], where rodents take alternating left and right turns ($left \rightarrow right \rightarrow left \rightarrow right \cdots$), in addition to spatial place cells, 'splitter cell' representations form, which fire preferentially on left/right trials. These are non-spatial latent state representations, since they disambiguate the same spatial location (central trunk) depending on whether it is predicting 'go left' or 'go right'. **(f)** The aliasing problem in graphs: by representing state via an observation, these two graphs would appear the same. **(g)** Sequence prediction tasks are sufficient to learn latent state representations, since identical observations can have different neighbours. Sensory sequences, and the associated actions, can come from both space and non-space (e.g. families). Some sensory predictions can only be done by knowing (generalising) certain rules e.g. `North + East + South + West = 0` or `Parent + Sibling + Niece = 0`.

vector calculations[40,41]) and naturally extend to new locations (albeit periodically). By a clever choice of *representation*, grid cells prevent the need for *computation*.

## Non-spatial state-spaces

While it is easy to intuit good state-spaces in physical space, this problem becomes less clear in non-space. One promising approach, derived from RL[42,43], is to cast spatial learning as understanding relationships on a graph (Figure 1b). In space, nodes of a graph define physical locations, and so edges between nodes exist if two locations are directly connected. Graphs, unlike literal maps, need not consider as-the-crow-flies distances, since roads and airplanes render distant locations more connected. Instead, they consider non-Euclidean distances based on the notion of 'connectedness'. This is a re-conceptualisation of a map in terms of its connections (topology), as opposed to distances (geometry). Significantly, graphs also formalise non-spatial problems (Figure 1b-c). Family trees, social networks, atoms in molecules, and many other problems, all consist of relationships between entities and can be represented with graphs. Nodes in the graph no longer represent physical locations, but instead non-spatial locations - for instance, Alice is Bob's grandparent in a family tree. See Box 1 for graph-based state-space models.

Graphs define state-spaces and so afford value-based RL. They also afford planning; starting with Bob (characterised by a vector $v$ - all zeros except a 1 at the Bob node element), and multiplying $v$ by $T$ ($Tv$; $T$ is the transition matrix, where $T_{ij}$ is the transition probability from state $j$ to $i$), gives a distribution over future states after one step. Similarly, multiplying again by $T$ ($T^2v$) gives the distribution after two steps. Repeating this process until a non-zero entry appears in the Alice node provides the shortest-path route between Bob and Alice (assuming you have the right transition matrix) - which is two since Alice is Bob's grandparent. Naturally, exactly the same tree search process works between two locations in space.

---

**BOX 1: RL STATE-SPACES, GRAPHS, AND GRAPH REPRESENTATIONS**

*The problem of building graphs for cognitive maps is the same problem as building state-spaces in reinforcement learning. Crucially, the state-space in RL is tightly linked to behaviour (through rewards, values and policies). However, once the state space is defined there is a further choice of how each state is actually **represented**. Clever choice of representation can reduce online value/policy computations. This has allowed normative mathematical theories to predict neural representations.*

Reinforcement learning is concerned with taking appropriate actions at specific states ($s$) to maximise the expected (discounted by $\gamma$) sum of future rewards $v(s) = \mathbb{E}\left[r(s) + \gamma r(s') + \gamma^2 r(s'') \cdots \right]$, where $s'$ and $s''$ are states following $s$. Bellman[33] realised that this is a recursive equation, since the right-hand side contains the left-hand side but one step in the future: $v(s) = r(s) + \gamma \sum_t P(s' \mid s, \pi)v(s')$, where $P(s_{t+1} \mid s_t, \pi)$ is the transition probabilities between states under a policy $\pi$. In essence, Bellman's equation says the value of the current state is the reward at that state plus the average value of states you can transition to. If you can *assign credit* to each state (like these equations do), then taking good actions is easy: just go to the neighbouring state with the highest value $v(s')$.

RL state-spaces define graphs with transition matrix elements $T_{ij} = P(s_j \mid s_i, \pi)$. One graph representation, the successor representation[44] (SR), is particularly relevant to cognitive maps[43,45]. The SR is a (discounted) sum of $n$-step transition matrices - $S = \sum_n \gamma^n T^n$. Elements of this matrix, $S_{ij}$, describe connectedness via all possible paths between two locations. Critically, if we represent connections between states in the world in terms of the SR-distance, then computing value is easy, since the SR is one half of the value computation[44] ($v = Sr$ where $v$ and $r$ are vectors whose elements are values and reward at each state).

Stachenfeld and colleagues[45] noticed that the rows of $S$ look like hippocampal place cells, and the eigenvectors of $S$ resemble entorhinal grid cells (Figure 2; similar to work demonstrating that the eigenvectors of place cell correlation matrices resemble grid cells[46]). Notably, SR makes many predictions about how both grid and place cells behave in different environments and tasks[45,47–49]. Critically, it also makes predictions of representations in non-spatial tasks[45,50,51]. Because it derives from a theory of learning, it can also account for behavioural phenomena that are otherwise hard to explain[52].

One prominent issue with SR, however, is its policy-dependence[53]. This means that when rewards move - or, worse, when obstacles appear - value calculations using SR are no longer optimal[53]. A recent model addresses this problem[54], using linear RL[55]. This model builds a representation for default behaviours that can be linearly updated when rewards change to approximate the new optimal policy. The required default representation (DR) resembles the SR, and can therefore be computed from grids cells. The model further provides a novel account of how to build world representations *compositionally* out of component cells representations (e.g. how grid and border cells interact to represent the insertion of a barrier)[56]. We return to this important issue in box 4 and related text.
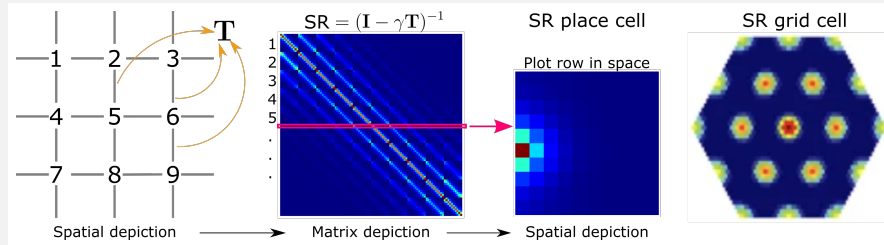
**Figure 2. Left:** Graphs can be succinctly represented via their transition matrix, $\boldsymbol{T}$. **Centre:** From this, the SR matrix can be calculated[44]. **Right:** Rows of this SR matrix resemble hippocampal place cells, and its eigenvectors resemble periodic entorhinal representations such as grid cells[45] (note this eigenvector is from a hexagonal, not a square, world).

## Latent states and sequence learning

Graphs are a flexible tool for representing problems, but how do we know which graphs to build? What should be the nodes in the graph or, equivalently, the states in the RL problem (Figure 1f)? A sensory observation alone cannot define a state, since two identical observations can have very different consequences; crossing a road implies looking right in the UK, but left in Germany. Formally, our world is not 'fully observable'; instead, we face 'partially observable' problems and must infer *latent state*[30,34] representations that disambiguate UK and German roads.

While individual observations are not enough to infer latent state representations, *sequences* of observations are, since all identical observations do not have identical surroundings (remembering you just read a newspaper in German is enough to tell you to look left when crossing a road (Figure 1g)). See Box 2 for clone-structured cognitive graph (CSCG) - a model that infers latent states from sensory sequences.

Neural representations in the hippocampal formation disambiguate states using latent representations[17,28,29,57–61]. For example, rodent grid cells will initially code identically for two identical boxes. However, after realising that the boxes are connected by a corridor, the grid representation changes to become consistent with the *global* two-box-and-corridor-space[28] (Figure 1d). These are latent state representations that disambiguate sensory aliased boxes due to their different futures. Physical location can also be aliased; in spatial alternation tasks[29,58] (Figure 1e), the same physical position (in this case, the central stem of the maze) predicts different futures depending on animal's left/right choice on the previous trial. Splitter cells[29,58], place cells[59], grid cells[28], lap cells[61], and others, are all cellular examples of the cognitive map disambiguating the world into *latent* states.

**BOX 2: BUILDING LATENT STATE REPRESENTATIONS FROM SEQUENCES**

*State-spaces must be inferred from observations. Because the sensory world is **aliased** - the same observation can occur more than once - states cannot be inferred from sensory appearance alone. Instead, **sequences** of observations uniquely identify states since two states with the same sensory observation will have different futures. States inferred via sequences are known as **latent** states, and building a latent state-space map can be used to afford different behaviours in sensorially identical situations.*

The clone-structured cognitive graph (CSCG) model[62] is an elegant approach for building de-aliased state-spaces. Here, hippocampus contains multiple 'clone' cells for each sensory observation[62,63]. Now, one hippocampal 'frog' clone cell responds to a frog in one location, and another responds if a frog appears elsewhere (Figure 3). The model uses Bayes to 1) infer which hippocampal clone cells should be active for each sensory observation and 2) learn an appropriate set of transition weights between clone cells. These transition weights are analogous to the transition matrix for graphs, but critically the state-space is *learned*, rather than provided by the modeller.

Many hippocampal findings can be understood in terms of representing latent states, from basic phenomena, such as place cells, through to complex representations which vary as a function of animal behaviour. These predictions are in the main common between CSCG and more complicated models that follow, and we show a number of these in detail in figure 6. A critical difference between CSCG and the following models, is that CSCG infers the whole latent space within the hippocampus (as opposed to the cortical input to hippocampus). This enables learning rules to be local, biologically plausible, and fast. By contrast, CSCG has to learn each map *de novo* and cannot benefit from having learnt similar maps before. It is exciting to think how these benefits may be combined (see section 'Complementary maps in hippocampus and cortex', Figure 7b).

CSCG is easily expressed in mathematics, and is closely related to hidden Markov models. From a sequence of

sensory observations $\mathbb{X} = \{\boldsymbol{x}_1, \boldsymbol{x}_2, \boldsymbol{x}_3, \cdots, \boldsymbol{x}_T\}$ and actions $\mathbb{A} = \{\boldsymbol{a}_1, \boldsymbol{a}_2, \boldsymbol{a}_3, \cdots, \boldsymbol{a}_T\}$, we wish to infer *discrete* latent states $\mathbb{Z} = \{\boldsymbol{z}_1, \boldsymbol{z}_2, \boldsymbol{z}_3, \cdots, \boldsymbol{z}_T\}$. Now, the same sensory observation, $\boldsymbol{x}$, can be linked to different latent states (clones) $\boldsymbol{z}$, via an 'emission' distribution $p(\boldsymbol{x} \mid \boldsymbol{z})$, naturally accounting for the aliasing problem. Along with predicting sensory observations, CSCG latent states predict future latent states and actions $p(\boldsymbol{z}_t, \boldsymbol{a}_t \mid \boldsymbol{z}_{t-1})$. Modelling the full sequence of observation is then:

$$p(\mathbb{X}, \mathbb{Z}, \mathbb{A}) = p(\boldsymbol{z}_0) \prod_t p(\boldsymbol{x}_t \mid \boldsymbol{z}_t) p(\boldsymbol{z}_t, \boldsymbol{a}_t \mid \boldsymbol{z}_{t-1}) \tag{1}$$

Here, each element of $\boldsymbol{z}$, $z_i$, is a 'clone' of a sensory observation (Figure 3, note we use $t$ for vectors in time and $i$ for indexing elements of each vector). Concretely, if there are 4 possible sensory observations, and 5 clones for each observations, there will be 20 elements to $\boldsymbol{z}$. The probability of observing a 'frog' given a 'frog clone' is defined as 1, but 0 given a 'snail clone' - $p(\boldsymbol{x} \mid z_i \in C(\boldsymbol{x})) = 1$ whereas $p(\boldsymbol{x} \mid z_i \notin C(\boldsymbol{x})) = 0$ if $C(\boldsymbol{x})$ are the clones of $\boldsymbol{x}$. CSCG marginalises over $\boldsymbol{z}$ and uses the expectation-maximisation algorithm to train the model[64] - that is, learn an appropriate set of transition probabilities $p(\boldsymbol{z}_t, \boldsymbol{a}_t \mid \boldsymbol{z}_{t-1})$ and infer $\boldsymbol{z}_t$. Once trained, this model can be used for planning by inferring a sequence of actions and observations conditioned on a start and end clone.
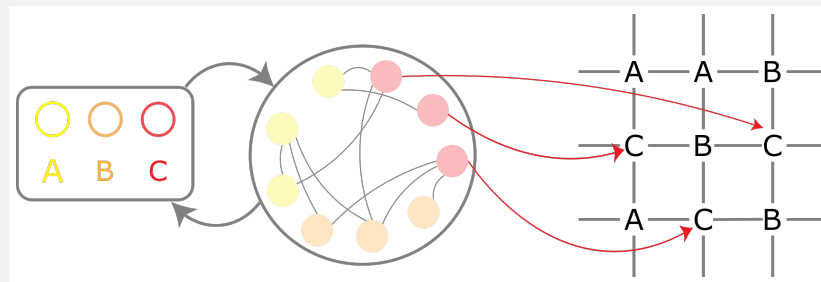


**Figure 3.** CSCG[62] addresses the problem of sensory-aliasing using multiple hippocampal clones for each sensory observation, then inferring which clone should be active at each location (as well as learning a clone transition matrix).

## Path integration and compression

Inferring latent states is really a problem of understanding where you are in an abstract space. In the two-room task[28], the global grid code which forms uniquely identifies physical locations. For the spatial alternation tasks, 'splitter' cells identify simultaneous location in physical space and location in left/right trials. Working out where you are in physical space is easy - accumulate self movement vectors (e.g. `North`, `South`, `East`, and `West` from head direction cells[65]) to update your relative position: this is path integration[66] (Figure 4a). Ants, rodents, birds, and humans all path integrate[67–69], with mammalian path integration crucially dependent on the hippocampal formation[70]. Entorhinal grid cells are considered an attractive substrate for path integration of two-dimensional spaces since 1) periodic representations extend to all, even unseen, space; 2) the periodicity of grid codes is inherently error-correcting[71]; 3) grid cells represent location with the same precision but far fewer cells than place cells[72]; 4) grid cells are experimentally driven more from path integration signals that place cells[73]. See Box 3 for path integrating models.

To work out where you are in graphs and non-space requires a modification of path integration. Rather than accumulating self-movement vectors, accumulate abstract movement vectors instead, (`Parent`, `Child`, `Sibling`, `Aunt`, `Nephew`, etc. for family tree graphs). Just like (x,y)-coordinates versus unique identifying representations for space, path integrating representations on graphs offer a benefit compared to representing every individual connection between nodes: adding a new node (Chloe is the `Sibling` of Bob) immediately implies all other connections (Chloe is the `Grandchild` of Alice) without needing to observe that relationship explicitly. This is because, just like (x,y)-coordinates, path integration treats all nodes equally and relationships are structured (`Sibling` + `Grandparent` = `Grandparent`). As such, only the few rules of path integration need to be known, not every possible relationship; path integration is a *compressed* representation. Not all graphs, however, can be path-integrated, since consistent actions do not always exist across graphs (for instance, social networks merely describe generic relationships).

*Path integration offers a powerful way to build latent state spaces. It builds maps that embed knowledge of the structure of the space (in physical space, `North` + `East` + `South` + `West` = 0). This means that path integration maps are: 1) inherently latent (and abstract), since they follow rules, not sensory observations and 2) allow relational knowledge to be transferred to any situations where the same rules apply. Notably, although path integration is not limited to space, not all graphs can be path-integrated.*

Path integrating models utilise a particular type of recurrent neural network (RNN) known as continuous attractor neural networks (CANNs[74]; Figure 4b), where neurons are recurrently connected via weights, $\boldsymbol{W}$, and receive velocity input, $\boldsymbol{a}$. The neural dynamics are given by:

$$\tau \frac{d\boldsymbol{g}}{dt} = -\boldsymbol{g} + f\left(\boldsymbol{W}\boldsymbol{g} + \boldsymbol{Ba}\right) \tag{2}$$

Where $\tau$ is the time constant of neuronal response and $f$ a non-linear activation function. With an appropriate set of weights, CANNs path integrate, with different cell classes (head-direction cells[74,75], place cells[76,77], grid cells[78]; Figure 4d) modelled with different weights. Remarkably, CANNs really exist in nature; ring attractors[79], both in connections and anatomy, are found in flies[80], and attractor manifolds are found in rodents[81,82].

Other path integrating models exist[83,84]. For example, velocity-coupled oscillators (VCOs) suggest path integration (along an axis) via interference between theta oscillations and velocity-dependent dendritic oscillations, with their phase difference indicating path integrated distance along an axis (this looks like a plane wave!). Here, grid cells are the sum of three such neurons with preferred axes at $\frac{\pi}{3}$ relative angles.

One major limitation of CANNs and VCOs, however, is that the weights of the recurrent weight matrix, $\boldsymbol{W}$, are carefully selected and not learned from sensory experience. However, it is easy enough to set up path integration as a learning problem via predicting observations $\boldsymbol{x}$: path integrate the latent state variable $\boldsymbol{z}$ and then predict observations $\boldsymbol{x}$ from the latent states:

$$p(\mathbb{X}, \mathbb{Z} \mid \mathbb{A}) = p(\boldsymbol{z}_0) \prod_t p(\boldsymbol{x}_t \mid \boldsymbol{z}_t) p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1}, \boldsymbol{a}_t) \tag{3}$$

Where the path integrating part ($p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1}, \boldsymbol{a}_t)$) is now replaced by a discrete-time version - i.e. $\boldsymbol{z}_t = f\left(\boldsymbol{W}\boldsymbol{z}_{t-1} + \boldsymbol{Ba}\right) + noise$. In fact, several models use a deterministic RNN (i.e. set the noise term to 0). These models successfully learn to path integrate when tasked with predicting ground truth spatial representations - i.e. $\boldsymbol{x}$ is either place cells[85], or (x,y) coordinates[86]. Neural units in both models form periodic representations (Figure 4e-f), but these are often amorphous, four-fold symmetric grids. An elegant analytic result[87], however, demonstrated that the 4- to 6-fold symmetry transition is governed by a single property: a third order regularisation term of grid cells. Indeed, this is easily implemented by the biological constraint of ensuring neural activity is positive[46,87].
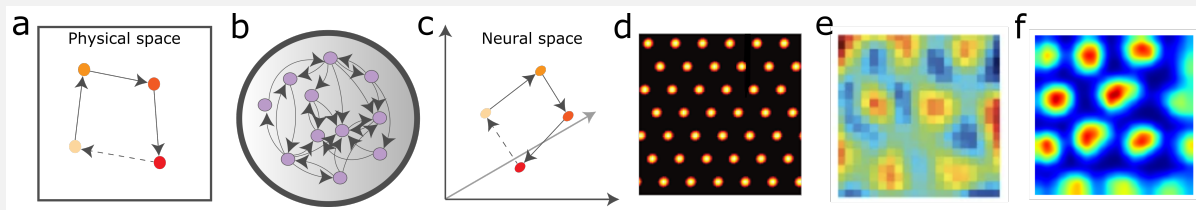


**Figure 4.** (a) Path integration is the summation of self-motion vectors to permit self-localisation. (b) CANNs are recurrently connected neurons which (c) path integrate, and (d) exhibit representations mirroring grid cells[78]. (e) In LSTMs[85] and (f) other RNNs[87], grid-like representations are learned when trained to predict ground-truth spatial representations.

---

An alternative, but less biologically plausible, equation is $\tau \frac{d\boldsymbol{g}}{dt} = -\boldsymbol{g} + f\left(\boldsymbol{W_a}\boldsymbol{g}\right)$, where the recurrent matrix $\boldsymbol{W_a}$ depends on the movement velocity.

## Generalisation

Generalisation, or the transfer of knowledge from one situation to another, is the substrate of the profound behavioural flexibility exhibited by animals. Without it, new situations could not be understood in the context of existing knowledge, and so prior

learned behaviours could not be leveraged in new situations. Generalisation in the sensory domain means realising that a Pekingese and a Rottweiler are both types of dog. In the structural domain, however, deep and powerful inferences can be made: doors often lead to new rooms; addition works for 100s as it does for 10s; the same path integration rules apply across different spaces. These pieces of knowledge have profound effects on behaviour.

Generalising with graphs, however, is hard as they require perfect alignment - if the bottom left of the new environment is prescribed the old environment's top right representation, then the rest of the environment cannot be represented, since the graph ends. Perfect alignment (*graph matching*, analogous to structural mapping[88]), however, is NP-hard and thus impractical in most situations. Generalising with periodic path integration representations, on the other hand, is easy since all positions are treated equally - representations corresponding to the bottom right in one environment could equally represent the middle of another environment. Furthermore since path integration maps are latent (thus abstract) by nature, they chart the relational structure of one family just as well as for another. This is generalisation of *relational* knowledge.

The hippocampal formation is critical for generalisation, along with memory, and some forms of imagination[1,2,5]. Hippocampal representations, however, do not generalise; neighbouring place fields are not necessarily neighbours in other environments, a phenomenon known as *remapping*[89–91] (Figure 1c). Entorhinal representations, on the other hand, do generalise; neighbouring grid cells (within-module) are still neighbouring grid cells in other environments - that is, the map is shifted and/or rotated (grid cell realignment[81,92]); representations corresponding to the bottom right in one environment now equally represent the middle of another environment. Spatial generalisation, at least, seems to exist in entorhinal cortex and is consistent with path integration.

Learning to generalise is often also a sequence-learning problem, but with sequences from *many different environments* (Figure 1g). When hearing about new family, after observing Daniel is Emily's parent, and Fran is Daniel's sibling, it is *only* possible to predict Fran's niece (Emily) with *a priori* learned relational knowledge: Parent + Sibling + Niece = 0.

To actually make sensory predictions, however, you need to know more than just abstract knowledge. You need to know how it interacts with real world representations (the abstract location in a family tree interacts (corresponds) with Emily for one family, and Chris for another; Figure 5a). One influential proposal is that hippocampal cells reflect this interaction, with abstract knowledge from medial entorhinal cortex (MEC) and sensory knowledge from lateral entorhinal cortex (LEC) combined rapidly ('fast-mapped'[93]) in hippocampus[9,27,94]. This bridges the abstract-to-real divide and permits generalisation, since the *same* abstract map can be reused across different sensory (LEC) environments and contexts. See Box 4 for models that *generalise*.

## Composition

Generalisation is more than transferring a single map to a new environment. Often only sub-components, or combinations of sub-components, need to be generalised. For example, understanding differently shaped rooms can be broken down into two components - an underlying 2D space and walls that can be placed anywhere. Should the cognitive map represent such common structural elements across tasks, then these elements can be *composed* to understand any given task configuration[27,54,56]. To encourage arbitrary composition, different structural elements (bases) should be represented independently (factorised) from one another[9]. Understanding a task then becomes a structural inference problem - finding the appropriate bases to represent the current task[95,96]; this is a form of structural analogy[97]. The value of such an approach has already been demonstrated in cognitive models which formalise composition in other domains of cognition, such as language and algebra[96,98].

The cognitive map of the hippocampal formation contains many such basis representations (Figure 1a,c). Object-vector cells[99] (OVCs), border-vector cells[100–104] (BVCs), reward cells[105], and goal-direction cells[106] (GVCs) , are all examples of *local bases* - representations that encode any object/border/goal, irrespective of where it is. Grid cells, by contrast, are examples of *global bases*, as they describe information equally across all space.

---

**BOX 4: GENERALISING WITH MEMORIES**

*We have seen models that build latent state representations, and models that path integrate. If these principles could be combined, we could build a powerful system that learns arbitrary latent states from sensory observations (like CSCG[62]) but additionally **generalises** these representations (like path integration models[78,85]) and composes them arbitrarily. For **abstract** representations to be reused (generalised) in different **sensory** environments, the **same** abstract locations must be 'linked' to **different** sensory observations. Hippocampal memories offer the ideal substrate for this link - they can rapidly tie sensory observations to specific locations.*

Hippocampal models of generalisation (the Tolman-Eichenbaum machine, TEM[27], and the spatial memory pipeline, SMP[107]) are tasked with predicting, as fast as possible, sensory observations in novel, but structurally similar, environments (for example, multiple different families or 2D worlds; Figure 1g). Both models consist of two key components: 1) An abstract path-integration module that is reusable across environments. 2) A *relational memory*[2] module that, like an address book, links abstract location representations with sensory representations (Figure 5a). These links change

from world to world, allowing the same abstractions to apply to multiple worlds.

Recall the probabilistic interpretation of path integration:

$$p(\mathbb{X}, \mathbb{Z} \mid \mathbb{A}) = p(\boldsymbol{z}_0) \prod_t p(\boldsymbol{x}_t \mid \boldsymbol{z}_t) p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1}, \boldsymbol{a}_t) \tag{4}$$

Previously, $p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1}, \boldsymbol{a}_t)$ was fixed and so each abstract location $\boldsymbol{z}$ could only predict a single sensory observation $\boldsymbol{x}$. If, instead, we had an address book of relational memories $\boldsymbol{M}$, we could remember what is where in *each* environment. To predict upcoming sensory observations, all that is required is to imagine a transition in abstract representation ($\boldsymbol{z}$, via path integration), then retrieve the memory at that location ('what' did I see the last time I was 'here'). Sensory prediction is now a combination of path integration and memory retrieval. But what space are we path integrating in, and how does it get built? Previously the weights, $\boldsymbol{W}$, in the path integrator ($p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1}, \boldsymbol{a}_t)$, where $\boldsymbol{z}_t = f(\boldsymbol{W}\boldsymbol{z}_{t-1} + \boldsymbol{B}\boldsymbol{a}) + noise$)) were built from predicting (x,y) coordinates or place cells (i.e. spatially curated representations). Now, we can predict actual sensory observations. This is more powerful. When sensory objects are arranged in space, the same spatial path integration as previous models will be learned, but when the sensory world has more complex dependencies, these will also be learnt. If the best way to predict the sensory future is to learn a complex map of latent states, then these models will learn to path integrate in this latent space (Figure 6).

While TEM and SMP are conceptually the same model, they have different implementations. Two critical ones are (1) that TEM is supplied with allocentric actions and object representations, but SMP must infer them from egocentric input and pixels; (2) that SMP implements memory with a memory network from machine learning[108], while TEM uses more biologically realistic Hebbian learning[109] and Hopfield networks[110]. This biological constraint means that the link between abstract and sensory world must take place in neuronal units. That is, the same hippocampal neurons must know *both* the abstract location and the sensory prediction. This type of *conjunctive* representation is commonly observed real in hippocampal neurons[59, 111].

TEM and SMP are deep artificial neural networks which learn to generalise structural knowledge and recapitulate a host of known representations of the hippocampal cognitive map in doing so (Figure 5c-d). Since SMP works from ego-centric inputs, it generates cells involved in the ego- to allo-centric coordinate transformation[101]. Additionally, TEM learns *compositional* entorhinal representations in spatial and non-spatial tasks, and can solve classical relational memory tasks that are crucially dependent on the hippocampal formation, such as transitive inference[57, 112].
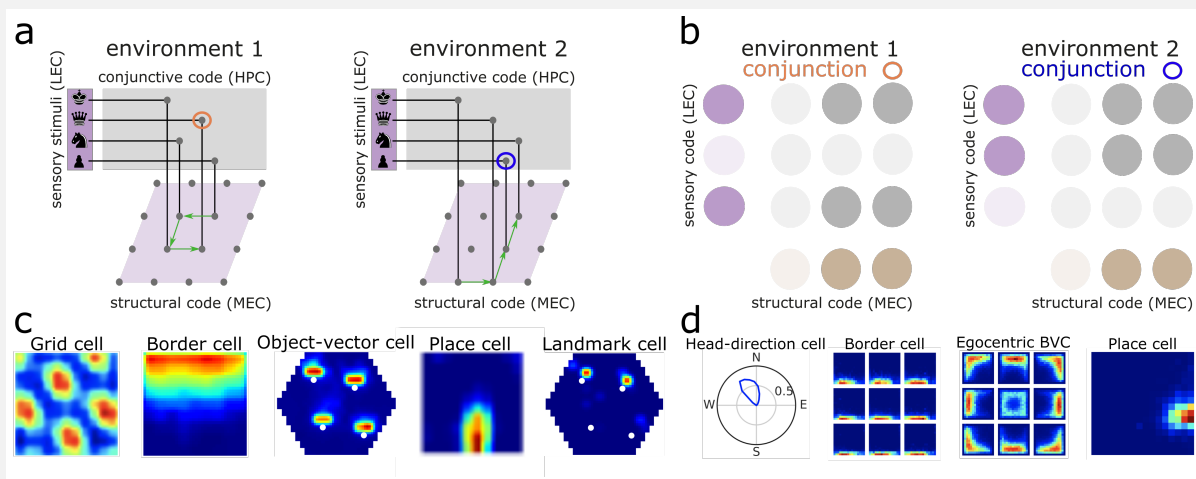


**Figure 5.** **(a)** Schematic of models (adapted from Sanders *et al.*[113]), where the *same* cortical representations (LEC and MEC) are reused in different environments, facilitated by *different* hippocampal combinations. **(b)** TEM conjunctive hippocampal cells receive input from particular MEC and LEC cells; hippocampal remapping affords generalisation. **(c)** TEM and **(d)** SMP recapitulate a host of empirically described cell representations[27, 107].

## Novel interpretations, integrations, and predictions

While the models above account for a variety of data, including hippocampal and entorhinal cellular representations from spatial and non-spatial tasks, they often do so in seemingly divergent ways, and there are many neural phenomena that remain

perplexing. Here we consider how these ideas can be integrated in order to model and understand cognitive maps at a deeper level, and offer novel accounts of several neural phenomena through a formal lens.

## Non-spatial hippocampal cells are latent state representations for generalisation

Many non-spatial hippocampal representations have appeared over the last few decades[17,29,58,61,114]. These cells appear to represent different tasks in different ways, but we can understand all these cell types with a single framework: representing latent state-spaces. We have argued latent state representations serve two purposes; firstly, separating states that have different futures and, secondly, affording generalisation, since the latent map can be used across multiple different environments. These arguments suggest two things: 1) Hippocampal and cortical representations do not map abstract spaces for the sake of it, but only to disambiguate states with different futures. 2) To generalise as fast as possible, every level of abstraction needs to be represented simultaneously; representing space in spatial tasks, non-space in non spatial tasks, and *both* space and non-space in interacting spatial-non-spatial tasks.

As a didactic example, consider spatial alternation tasks[29,58] where animals cycle $left \rightarrow right \rightarrow left \rightarrow right \cdots$ at a choice point (Figure 6a). This task can be 'un-rolled' into a 'big-loop' state-space where the first half is going left and the second is going right. This is a latent state-space for the task; it de-aliases the common 'trunk' section according to whether the animal is going to take a left or right turn. This 'big-loop', however, ignores spatial knowledge - understanding the big-loop *alone* does not let you *know* you are back in the same place each time you return to the common section - to generalise spatial knowledge you additionally need a spatial representation. Hippocampal cells in this task indeed code for both space (place cells) and big-loop (splitter cells)[29].

This interpretation of non-spatial hippocampal representations (latent states for disambiguation and generalisation) can be modelled formally. Here we show many non-spatial tasks[17,29,58,61,114] can be understood by these two principles alone. We use TEM, since it learns and generalises latent states at multiple levels of abstraction. We note that pure latent state models, such as CSCG, account for the separation of states with different futures (e.g. it learns splitter cells for the big-loop), but without extra assumptions it will not learn place cells at the same time, as it cannot profit from generalising the structure of space.

Firstly, training TEM on spatial alternation tasks[29,58,114] (Figure 6a-c), TEM recapitulates both splitter cells and place cells. Splitter cells for the 'big-loop' and place cells for spatial generalisation. Secondly, when rodents are trained on a task where reward is received every 4 laps of a loop, hippocampus[61] contains both spatial 'place' cells that care about location in the lap, and non-spatial cells that additionally care about *which* lap. This is exactly what TEM learns, too (Figure 6d) - lap cells for the 'big-loop' and place cells for spatial generalisation. Lastly, when animals are trained to make left/right choices on a T-maze depending on the difference in number of sensory cues appearing on the left/right as the animal moves forwards in the central trunk (Figure 6e), hippocampal cells form an abstract map spanned by physical space and cue difference (termed 'evidence'). TEM learns exactly this; physical space to predict when to make the choice, and cue difference to predict reward left or reward right.

## Complementary maps in hippocampus and cortex

The reviewed models mirror an age-old debate in the hippocampal literature - does hippocampus map space[3], or is its role one of memory[1,115]. TEM and SMP suggest memories since they adhere to hippocampal indexing theory[116], where the hippocampal representations are simply an index that binds together cortical representations (structural (MEC) and sensory (LEC) representations in the case of TEM and SMP). In these models, hippocampus acts as a 'memory map', forming memories at distinct locations in the map, but the notion of 'map' is inherited, with all predictive capabilities occurring through cortex. These models embody an extreme version of Eichenbaum's view that the hippocampal role in navigation is principally one of relational memories[115]. By contrast, in the state-space models (such as SR, DR, and CSCG), the hippocampus is an explicit map, where connections between hippocampal cells define the connections of the map. Regardless, both sets of models explain many hippocampal representations and phenomena.

The observation that the models of generalisation (TEM and SMP[27,107]) treat hippocampus as memories alone, while models of single environments (SR, CSCG, DR[45,54,62]) treat hippocampus as a map, is a distinction that offers a potential unification of the hippocampus role in mapping and memories: it is easier to learn how to generalise if each (latent) state-space is already built. More precisely, should all states of the world be appropriately separated, and relationships between states known, cortex can receive high-fidelity training signals (since predictions can be compared to a de-aliased state-space), thereby significantly reducing the burden of learning. This means entirely novel sets of relationships can be efficiently learned as follows: first, form temporary hippocampal maps; next, learn the statistical structure of these maps in cortex (Figure 7a). This proposal follows complementary learning systems theory[117,118], where cortex slowly learns the statistics of hippocampal episodes. We take note of an interesting model[119] that, while not involving structural learning or generalisation, leverages two independent systems for self-localisation: relational maps in hippocampus and path integrating maps in entorhinal cortex.

This integrated approach is realisable within the existing models (Figure 7). CSCG can disambiguate sensory states and
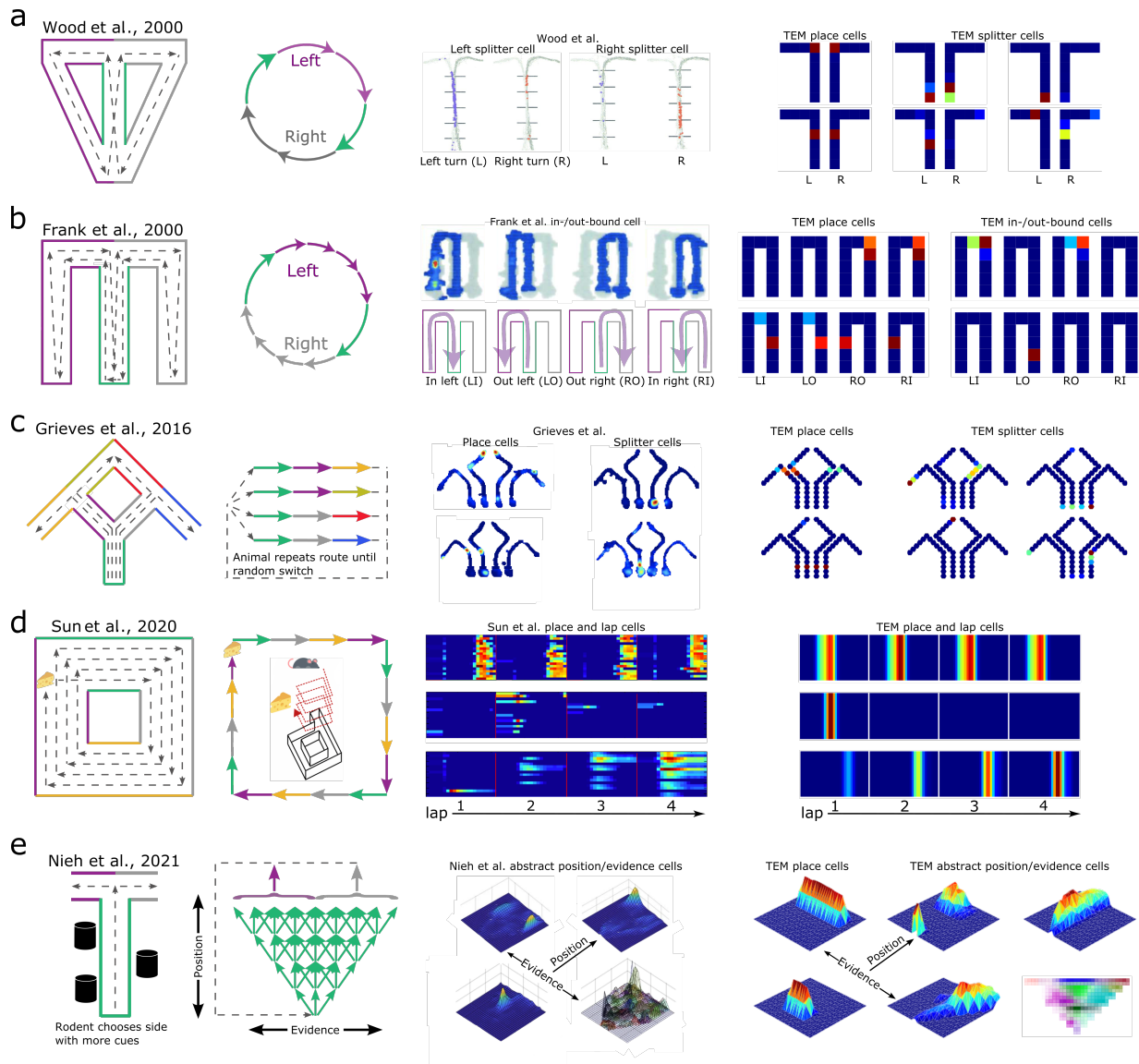
**Figure 6. Representing latent states.** Many apparently different neural phenomena are captured with a unifying computational principle - building state-spaces that can accurately predict different futures (latent states) as fast as possible (generalisation). **(a-e)** For each row, left/center-left are the task and its latent state-space (with colours denoting sensory experience), while center-right/right are real/TEM neural representations. **(a)** In a T-maze task[29], where animals alternate left/right turns, the state-space is described by a 'big-loop' latent space, since the central trunk predicts different futures depending on a previous left/right turn. Hippocampal cells represent both space (place cells) and the 'big-loop' (splitter cells). Splitter cells are trajectory-dependent, firing at the *same* spatial location (central trunk) *differently* depending on the prospective future (left/right). TEM learns both spatial (place) and non-spatial (splitter) cells when trained on this task; splitter cells to represent latent state in the 'big-loop' and place cells to represent physical location, and thus facilitate spatial generalisation. **(b-c)** More complicated spatial alternation tasks[58,60] are also described with 'big-loop' latent state-spaces. Both real and TEM hippocampal representations contain spatial (place) and non-spatial (trajectory-dependent) cell representations. **(d)** Performing 4 laps to receive a reward is a non-spatial task[61]. It is also described as a 'big-loop' latent state-space. Rodent hippocampus, and TEM, represent both space (place cells; top) and non-space (lap-specific cells; middle/bottom). **(e)** A T-maze task where rodents choose left/right depending on sensory evidence (as the animal moves along the central trunk) has a latent state-space spanned by position and evidence. Hippocampal cells, and TEM learned hippocampal representations, map this position-evidence latent space i.e. not just spatial location (the bottom right panel is a collection of many different cells representations). We note that CSCG would learn the non-spatial cells (e.g. splitter cells) but not the spatial cells, as it is not a model of generalisation. Code for simulations will be made available on publication.

rapidly learn hippocampal relational maps for *single* tasks, while TEM (or SMP) learns path integrable abstractions that can be generalised between *many* structurally similar tasks. Since both TEM and CSCG utilise multiple 'clone' hippocampal cells for each sensory observation, it is particularly easy to combine these models. This would be formulated as a TEM-like model, but where hippocampus is predictive of future hippocampal states (Figure 7b). Such an approach combines the best of both models - learning novel maps fast (CSCG), but also leveraging past knowledge to understand similarities between maps (TEM/SMP).

## Cognitive maps and behaviour

The models discussed here interact with behaviour in different ways (see Box 5 for additional discussion on using eigenspaces for various behaviours). The models formulated in the RL framework (SR and DR[45,54]) can be explicitly used for model-free and model-based learning, and so provide insight into the hippocampal role in constructing state-spaces for RL; constructing a *predictive* cognitive map[43,45] [ii]. Since these models only require well-separated state-space as input, any model that learns and builds state-spaces, such as CSCG, SMP, and TEM, could also act as input to these RL models. For example, the discrete state-space of CSCG (or TEM hippocampal cells) can be used for SR (Figure 7c).

The sequential models, such as CSCG, offer an alternative to tree search - 'planning by inference'[39,121]. This involves conditioning a probabilistic model on a start and goal state, then inferring a distribution over action sequences and their intermediary states. Since CSCG is a Bayesian model, it can naturally implement this procedure, thereby suggesting a hippocampal role in action inference. In principle, any probabilistic model with states and actions can be used to perform action inference[122] (e.g. TEM), but the practicality of this differs across models and the details of their implementations.

For models that learn grid codes, vector-based planning can be used (for example, SMP and other neural network models implement vector-based navigation[40,85,107]). These models perform navigation and short-cutting behaviours reminiscent of biological agents, along with transferring policies from one environment to another (which is possible because these representations generalise).

---

**BOX 5: EIGEN-SPACES**

*The observation that grid cells resemble eigenvectors of place cells (or of the spatial transition matrix) has led to interesting suggestions about mechanisms for planning and exploration.*

To plan the future, you need to look across multiple transitions. Eigenvectors simplify this problem because all multi-step transition matrices share the same Eigenvectors. We can see this by diagonalising the transition matrix $\boldsymbol{T} = \boldsymbol{V}\Lambda\boldsymbol{V}^T$, where $\boldsymbol{V}$ is a matrix with eigenvectors as columns, and $\Lambda$ is a diagonal matrix of eigenvalues. A 1-step transition of a state vector $\boldsymbol{s}$ ( $\boldsymbol{T}\boldsymbol{s} = \boldsymbol{V}\Lambda\boldsymbol{V}^T\boldsymbol{s}$ ) uses the same eigenvectors as a 2-step transition ($\boldsymbol{T}^2\boldsymbol{s} = \boldsymbol{V}\Lambda\boldsymbol{V}^T\boldsymbol{V}\Lambda^n\boldsymbol{V}^T\boldsymbol{s} = \boldsymbol{V}\Lambda^2\boldsymbol{V}^T\boldsymbol{s}$ since $\boldsymbol{V}^T\boldsymbol{V} = \boldsymbol{I}$), and so on. The eigenvectors for *local* and *non-local* transitions are the same.

Intuitively, this means these eigenvectors can be used for exploration, planning, sampling in replay, or any other type of multi-step navigation. Different sampling patterns differ only in the eigenvalue matrix $\Lambda$. When both eigenvalues and eigenvectors come from simple diffusion, navigation/ exploration strategies visit surrounding locations one by one, progressing further afield slowly just like diffusion. On the other hand, by making a clever alternate choice of eigenvalue matrix (using a bespoke diagonal matrix, $\Upsilon$, rather than the diffusion eigenvalues, $\Lambda$, but still using the same eigenvectors; $\boldsymbol{V}\Upsilon\boldsymbol{V}^T\boldsymbol{s}$ ), very different strategies emerge[123], such as turbulence or super-diffusion (also known as Lévy flights, which are known to be used by animals exploring environments for food[124]), and can be seen in rodent hippocampal replay[125].

In fact, with another choice of weighting matrix, $\Upsilon = \sum_n \gamma_n \Lambda^n$, you can exactly compute the SR under a diffusive policy, which is closely related to the distance between states. This is particularly interesting, as when you have distances planning is easy - just go the the neighbouring state with lowest distance from the goal. Importantly, this planning is by eigenvectors (grid cells) alone - in fact all you need for this 'intuitive planning'[126] are the start and goal grid codes.

So far we have been considering diffusive transition matrices, i.e. matrices without actions. However, by making transition matrices actions dependent (remember path integration has action dependent matrices too) we can play games just like path integration. In space, at least, the transition matrices needed for different actions all have exactly the same eigenvectors, but different (complex) eigenvalues. Hence, path integration can be reduced to successively adding the eigenvalues associated with each action[127]. This way of thinking unifies path integration with SR-like planning. Interestingly, it also brings different models of path integration into a common framework since, in this case, the eigenvectors are plane waves (not grids as the transitions are unidirectional!) just like those required for VCOs[83,84], and the transition matrix is just like the weight matrices required for CANNs[128].

---

[ii]It remains unclear, however, whether hippocampal neurons could represent the vanishingly small differences in SR between (sometimes adjacent) states necessary for accurate reward-guided behaviour.
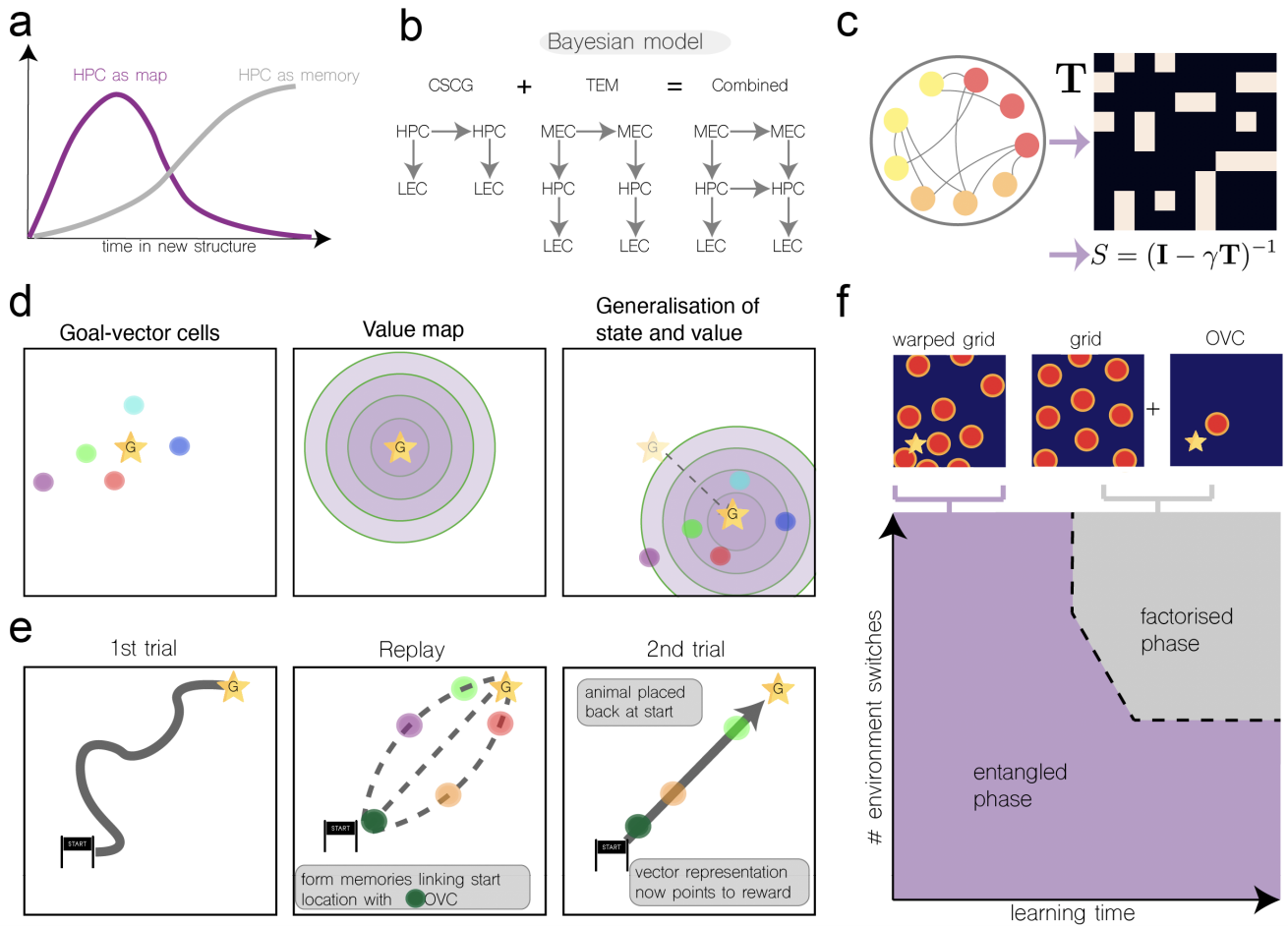
**Figure 7. Integrating different cognitive map models and novel predictions. (a-b)** The reviewed models suggest two roles for hippocampus: 1) A map - connections between hippocampal cells encode relationship between states - and 2) memories linking *cortical* map representations. **(a)** We suggest hippocampus serves *both* roles, but does so in different situations. In experiences where no prior cortical map is useful, hippocampal representations build a relational map; in familiar experiences where cortex has already learned how to structure (e.g. path integration) representations, hippocampus fulfils the role of memory. We suggest that with increasing experience, there is a transition from hippocampus as a map to memory, and this will be tied to the behavioural ability to generalise (via cortex). **(b)** TEM (HPC as memories) and CSCG (HPC as map) models can be easily integrated (since both are formalised probabilistically) into a model with a hippocampus that can form both maps and memories. **(c-e)** State-spaces for behaviour. **(c)** Learned latent state-spaces can be inputted into RL algorithms such as the successor representation. **(d)** On the other hand, compositional representations, such as GVCs, permit rapid generalisation of policy. Since these representations already generalise to novel goals in novel environments, all that is required is a pre-computed set of values (or policies) associated with the GVCs. The value map (or policy) is simply transferred along with the GVCs: *credit assignment through generalisation*. **(e)** Replay might play a role in this mechanism. After encountering a goal, we want the goal-vector representations to exist across all of space, and especially any start locations. Replay trajectories provide an offline solution; path integrate (offline) GVCs and bind them (via memory) to important locations such as the start state. Thus, when re-entering the same environment, vector representations and the associated value map (or policy) already exist. This is replay as the offline building of maps for credit assignment through generalisation. **(f)** The aforementioned mechanisms rely on compositional (factorised) representations. Sometimes, however, brain representations are not compositional, but entangled[120]. Since compositional representations are beneficial for generalisation, we suggest animals have factorised or entangled representations depending on the pressure to generalise; regularly staying in the same task will encourage entangled representations, while regularly switching tasks will encourage factorised representations.

## Credit assignment through generalisation and the interplay with striatal RL

Credit assignment is the attribution of value to state. RL typically assumes that the underlying state-space is fixed, and values are *slowly* assigned to these states. There is no requirement for state representations to be fixed, however; they can change to better represent value. For example, after encountering a goal, GVCs (goal-vector cells) form[106] - cells that are active at certain distances and directions from goals (Figure 7d). This can be interpreted as a state representation augmentation (with new cells). Importantly, since GVCs path integrate, once a single GVC forms at a goal, all others can be built for free as the animal navigates the map.

We propose that such compositional representations come 'pre-credit assigned', such that they can *immediately* provide a state-space that accurately predicts value (or policy). This can easily be understood in the context of tasks with changing goal locations; pre-learned goal-vector representations can be immediately composed with spatial representations to generate an accurate and flexible representation of any goal state (Figure 7d), and since these GVCs come with value (or policy) already attached, optimal actions can be taken much more rapidly than in the case of traditional RL-learning [iii]. The only *online* role of the cognitive map is inferring which pre-learned and pre-credit assigned representations to compose. In the case of goals this is easy - only GVCs are required, though it can be other *local bases* in other circumstances. This is **credit assignment through generalisation**, and is akin to meta-RL[129, 130], since prior statistical knowledge (e.g. GVCs) can be integrated on-the-fly to solve novel tasks.

Where do these representations come from in the first place? The cognitive map models suggest that such representations can be learned from statistics of behaviour. Just as OVCs can be learned as reflections of regular movements towards objects[27], GVCs can be learned when behaviour is biased towards goals. In general, to train these 'pre-credit assigned' compositional representations, cortex must learn from sequences of behaviour. This suggests an interesting interplay between generalisation and reinforcement learning. Initially, behaviour is generated via classic RL (perhaps in the striatum). Understandably, initial striatal actions will be bad (when encountering entirely novel tasks), but as RL learns good policies, actions will be towards goals. The cortico-hippocampal system can then learn compositional representations of these policies (e.g. GVCs) from the statistics of these sequences. In novel tasks, behaviour can be generated entirely from generalisation, with no need for new striatal RL. This also offers a virtuous cycle, where learned general cortical representations can be provided back to striatum as a state-space for RL, and so on. The notion of offline cortical learning from striatal actions sequences relates to recent machine learning methods in offline RL. Here, sequence models learn the statistics of behavioural sequences from conventional RL algorithms, after which the sequence model can be used for planning[131, 132] in a manner analogous to planning by inference. In sum, this proposal offers a novel role of cortical-basal ganglia interaction for constructing RL state-spaces and generalising policies.

### Replay: offline state-space construction

If behavioural control in a new world is reduced to a state-space composition problem, it becomes important to construct state-spaces rapidly and accurately, and to store them in memory so they can inform future decisions. To build such memories requires path integration (for example, to ensure the correct goal-vector cell is tied (composed) to the correct location), but to build them quickly this composition should be done as much as possible offline - it should not be the animal's actual location that is integrated.

An appealing substrate for this composition is replay[133]. For example, when an animal receives reward, it is important that all other states in the environment are aware of their relative location to the reward. Replay can path integrate away from the reward, successively tying (composing) each new goal-vector cell to its respective hippocampal/cortical location (perhaps building landmark cells in hippocampus[134]; this is a similar mechanism to the simultaneous grid and place cell replay from Evans and Burgess[119], but now used to instantiate rewarding policies, instead of ensuring consistency between place and grid representations). Now, should the animal return to a state, that state representation already 'knows' about its relation to the reward (Figure 7e). It is no longer necessary to hold all goal locations in mind, as the state-space composition is stored in memory. All heavy computations of building state-spaces take place off-line, and thus the computational burden is reduced for online behaviour. This idea relates to previous ideas from RL that cast replay as a mechanism for optimal credit assignment to existing states[135], or a mechanism for building state-spaces from scratch[32, 52]. However, in a generalisation framework (outlined in the section above), these two computational processes are subsumed by the single process of composing state-spaces from pre-learnt bases. To test this framework against data, it will be interesting to build a formal understanding of optimal replay patterns under these assumptions. Notably, it will make predictions not only about patterns of hippocampal replay, but also if and when these patterns will align with replay of more abstract representations in entorhinal and frontal cortices[136, 137].

---

[iii]This relates to DR, but here representations are truly compositional - the same GVCs can be used anywhere in space! Additionally, values can be optimal and not just optimal under linear RL assumptions.

## When neural representations factorise

Grid cells were once thought of as representations of space and space alone. By apparently ignoring sensory (or other non-spatial) details of an environment, grid cells were considered a *factorised* representation of space. Similarly, other spatial representations found in entorhinal cortex, such as OVCs and BVCs, are seemingly factorised, since they compositionally augment the entorhinal grid representation to represent different environment configurations. Recent evidence, however, has shown that grid cells warp towards consistently rewarded locations[120, 138]. Factorised representations do not warp, since warping is an environment-specific phenomenon; warping around rewards does not transfer to different spatial configurations of rewards.

Why should grid cells warp and sometimes not? There is a computational trade-off between using factorised compositional bases and using bespoke warped representations - specifically, a pressure to generalise versus precisely representing a single task (Figure 7f). With infrequent task switches (that is, repetitively solving a single task), it is more efficient to learn and store a bespoke warped representation (warped since the animal's notion of space becomes warped as it just does a stereotyped looping behaviour in space), as generalisation is not necessary and storing one representation is more efficient than combining many. With regular task switches (such as solving different goal configurations of the same task), the pressure for generalisation is high, and so *compositional* bases are favourable. This idea can be stated succinctly: when the set of tasks that an animal faces is itself factorised, then cellular representations for that task will also be factorised (made up of compositional bases) (Figure 7f). This hypothesis can thus make simple and falsifiable predictions in spatial tasks with environmental rewards. When rewards and space regularly occur in any combination (factorised), both representations of space (grid cells) and reward (reward-vector cells) will exist. By contrast, when rewards and space always occur in the same combination, a bespoke, warped representation will suffice[120, 138].

# Open questions

## The role of time in memory and cognitive maps

The discussion of cognitive map models so far assumes that learned representations remain stable over time. This clearly cannot be the case, since we can remember events at the same place and same conditions but on different days - hippocampal memory representations (e.g. place cells) must be different for different moments in time. An accumulating body of evidence indicates exactly this, with neural representations *drifting* over time and experience (Figure 8c), challenging traditional notions of engrams and receptive fields[139–142].

But how can hippocampus maintain a stable representation of space, if the cellular basis of this representation is drifting over time? Generalisation models offer a natural solution as, here, hippocampal cells bind multiple factors of the input. Only one factor needs to change for the entire hippocampal representation to change (Figure 8d). If entorhinal cortex could learn abstracted representations of time as well as space, then as time passes the temporal code will progress and the hippocampal code will drift to new cells, but these new cells will only differ in their connections to the entorhinal cells that represent time, not those that represent space (Figure 8d). Representational drift, in this view, is just hippocampal remapping, but now it is not sensory observations or space that has changed (as in Figure 5b), but time instead. A prediction that follows is that the order of drifting cells is not random.

Hippocampal represents time thought more than just drift. Pure 'time cells', for example emerge when rodents are required to stay still, or run on a wheel, for a particular duration of time in a task (Figure 8a)[143, 144]. These cells can be easily understood as enabling prediction of when the delay period finishes (Figure 8b). Crucially, this temporal representation is just one part of a overall map relating experiences to one another. More precisely, during the delay period, while space is not changing, *position in task* is changing. It is this overall task position that cognitive map models suggest is being represented in 'time cells'. Indeed, as we have demonstrated in Figure 6, neurons representing latent states often appear to be tracking *progress through task*, and in this view, time cells can be interpreted as another instance of this.

Lastly, viewing representational drift through the lens of temporal abstraction might help cognitive maps build abstractions from limited data. In particular, since learning by extracting the statistics of *multiple* different situations is critical for abstraction, drifting representations may allow a *single* environment to instantiate *multiple* representations, appearing (to cortex) as if they were from multiple different environments (Figure 8d). This process relates to data-augmentation from machine learning[145], in which additionally presenting various transformations (such as rotation, cropping, colour distortions, noise addition) of the *same* data dramatically improves representation learning and generalisation, since the abstract quantity of interest is invariant to the various transformations (a rotated dog is still a dog, for example). Similarly, should the structures of interest (within the sequential input data - e.g. loops, space, task) be invariant to drifting hippocampal representations, then abstraction and generalisation could be enhanced. In sum, casting representational drift as a form of *temporal abstraction* for generalisation might prove fruitful, and integrate the previously confusing findings of representational drift with existing models of cognitive maps.
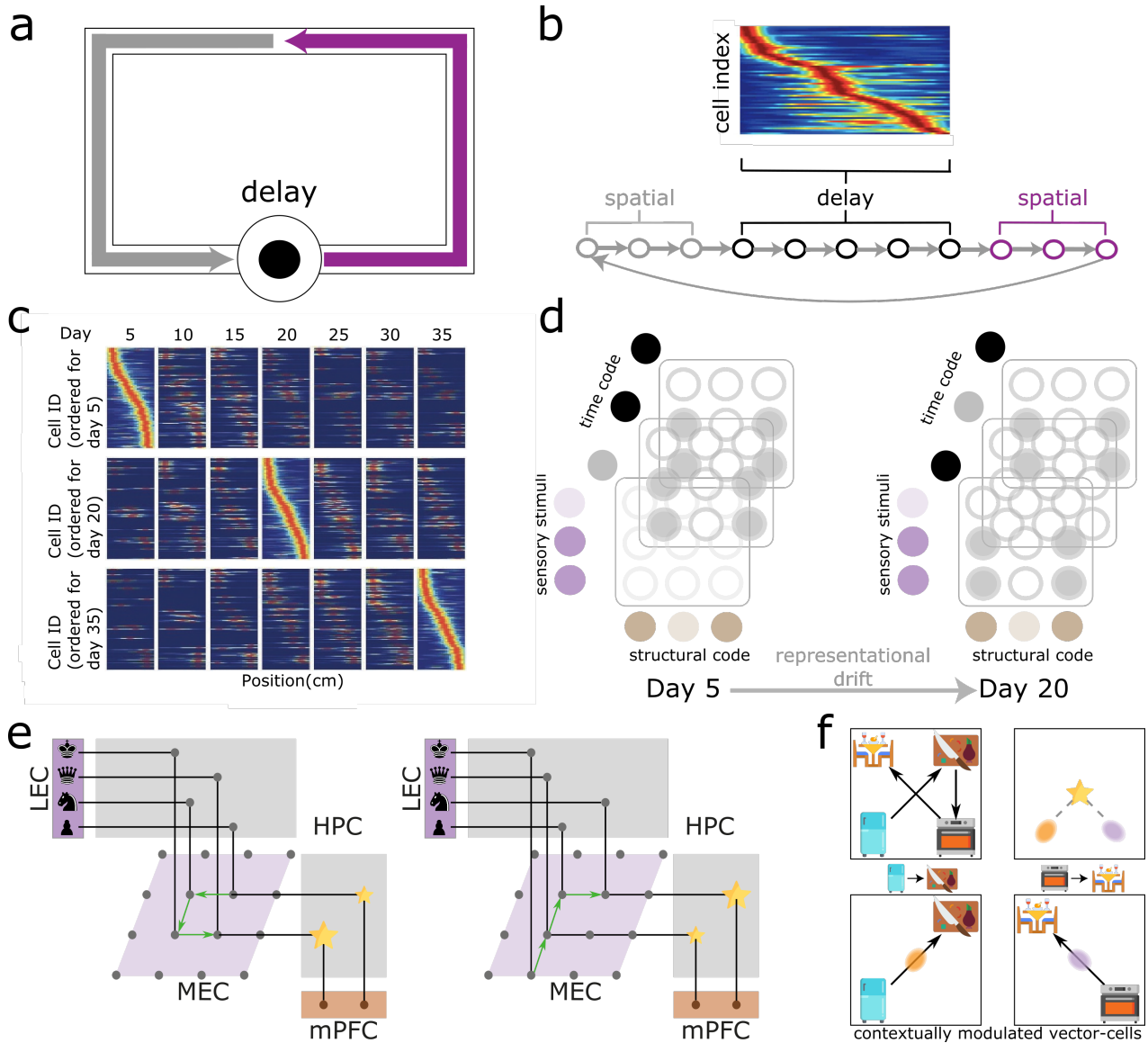
**Figure 8. Representing time and hierarchies of abstraction in cognitive maps. (a-b)** Neuronal representations of time might structure tasks according to 'progress' through that task. **(a)** In a task with a delay period, the full **(b)** latent state-space of the task includes the delay period, since 'progress' through the delay must be represented to predict when the delay period ends: in essence, latent states for predicting the future. Indeed sequences of hippocampal cells fire during the delay period as if they were coding time[143,144] (adapted from Salz *et al.*[146]). **(c)** Time additionally impacts representations via drift[139]. Here, hippocampal (and other) representations slowly change over many days, such that an entirely different representation encodes the same location (adapted from Ziv *et al.*[139]). **(d)** While the mechanism, or function, of drift is unknown, a tantalising possibility, inspired by the reviewed models, is that representational drift is remapping in disguise. In particular, as in TEM, if hippocampal representations reflect a triple conjunction of space, sensory stimuli, and time, then drifting hippocampal representations can parsimoniously be due to changing time representation while space and sensory representations remain the same. Rather than spatial remapping (Figure 5b), this is *temporal* remapping. **(e-f)** Representing hierarchical tasks. **(e)** Schematic of a *hierarchical* version of TEM, where an additional prefrontal module is included. Should this module represent location in task at an abstracted level (e.g. 'just before the over' in a recipe), this abstract and *non-spatial* representation can contextualise the hippocampal-entorhinal system - set goals *in space*. We note that this schematic is a simplification of true neuroanatomy (e.g. mPFC-HC connections may go via reuniens). **(f)** This predicts novel cell types, such as route-dependent GVCs: representations that point towards goal locations but only at specific points in the task (e.g. only before chopping the vegetables). This is analogous to splitter cells, though these representations can occur anywhere in space, not just at specific points on a T-maze. Icons from https://www.flaticon.com.

## Interacting levels of abstraction

We have shown how models can build abstract representations that generalise over different *sensory* realisations, but the real power of abstractions comes when this process can happen repeatedly, so that abstractions can themselves lead to further abstractions. When we are learning to cook a new recipe, we don't need to relearn the rules of space to find the oven, and when the recipe is learnt it can easily be transferred to kitchens with new spatial layouts.

In the latent state tasks from earlier (e.g. spatial alternation; Figure 6), while there was both a task ('go left then right') and space at play, they came in a fixed configuration; the latent space would not have generalised if the T-maze became a W-maze. Cooking recipes in different kitchens is equivalent to a T- to W-maze switch, thus we need something fundamentally new in the models to account for this. One attractive options is for *both* spatial representation and task ($left \rightarrow right \rightarrow left \rightarrow right \cdots$) representations need to be separately represented (factorised) so they can be arbitrarily combined (ovens being in different locations in different kitchens). Given enough experiences of different kitchens, this factorisation can emerge from training, and does in many deep learning networks[147]. However, using the same tricks as before - i.e. using the hippocampus as a mediator of factorised representations (e.g. space and task) - the required number of recipes and kitchens for training can be dramatically reduced. One intriguing possibility is that the different representations observed in fronto-temporal cortices[148–153] might reflect such a factorisation, with entorhinal representations grounded in interactions with the physical environment, while neurons in PFC representing abstract, task-related invariances, such as 'location in task'[35, 137, 148, 153–155].

While a factorisation allows representation of any space-task combination, to actually understand any given space-task combination, these representations must interact. The go-to-oven mPFC representations needs to be linked to the spatial location of the oven, or vector cells pointing towards the oven, in order to navigate to the oven. This linking can occur in exactly the same way as the earlier models suggest - through hippocampal memories (Figure 8e; we note the mPFC-HPC connection is likely mediated, for example via nucleus reuniens[156]). Interestingly, though, the very same vector cells can be reused whether it be the oven or the chopping board. This make a prediction - vector cells that are contextually modulated[150] depending on 'location in task' (Figure 8f).

Building models of interacting task and spatial representations, with principles of abstraction, generalisation, and path integration, allows neural representations from RL tasks to be understood in the *same* language as space. With emergent task-level (mPFC) representations potentially revealing insights into how cells might represent task structure itself, they will therefore be of interest whenever animals are shaped to perform tasks.

## From sequences to other domains of cognition

The models we have described translate the problem of building maps into problems of understanding the structure of possible sequences. This raises two interesting points. Firstly, there are many other sequence problems we face that are not traditionally thought of as related to cognitive maps - perhaps these can be understood similarly to space and tasks[157, 158]. Secondly, sequence problems are not the only cognitive problems we face - how can the understanding we have got form sequences extend to these domains?

How far can we get with sequence problems? Machine learning has taught us that sequence learners (recurrent neural networks, long short-term memory units[159], Transformers[160]) can perform well on a wide variety of tasks including language processing, mathematical understanding, and logic problems[161, 162]. This makes sense since language, mathematics, and formal logic are *sequence* problems where generalisation is key[163, 164]. Each comprises of content (words/numbers) combined within different structures (grammatical rules/ mathematical operators), and vice versa. Whilst mathematics and language engage large (and different) cortical territories[165], it is interesting to consider whether the neuronal representations that support these functions might be understood with principles similar to state representation, factorisation, and path integration described in the sections above. For example, mathematical operators, like addition and subtraction, bear similarity to forwards and backwards actions on a line (similar analogies can be made for integration and differentiation). One intriguing finding is that hippocampal and entorhinal cells have fields that respond to unique numbers[166].

What about non-sequence problems? Much of the neural processing underlying cognitive problems does not seemingly require sequence transitions. For example, understanding a football and the Earth are both spheres does not require learning from sequences; thus, it is not clear if organising principles similar to space play a role in learning these abstractions. Analogies between path integration and understanding spheres, however, can be made. The data-generative factors of a ball - {sizes, shapes, colours} - are all examples of variables which can be projected onto a manifold where 'actions' such as `add-red`, `bigger`, `remove-red`, then `smaller` have a meaning (and would take you back to the same sphere!). Indeed modern machine learning methods learn such manifolds from images inputted in no particular sequence[147, 167]. Some non-sequential problems can also be reformulated sequentially. While an image is not sequential itself, it can be viewed sequentially. In this vein, it is notable that grid-like cells have been observed in monkey[168] and human[169, 170] entorhinal cortex during saccades on images. Similarly, when humans view silhouettes of stacked objects, component objects are replayed sequentially[171].

## Conclusion

The hippocampal formation is a poster child for cognitive neuroscience because of its beautifully organised neuronal responses and the profound effects of its damage. However, whilst these experimental findings seem self-explanatory when examined in simple situations like open-field foraging, they have been hard to relate to complex real-world behaviours. Moreover, it has not been clear whether the amazing discoveries gleaned from examining rodents navigating in space might have broader implications for understanding more general cognitive processes. The ideas reviewed in this paper offer concrete and formal methods for addressing these long-standing questions. Excitingly they do this by re-imagining the problem. By asking questions such as 'what really is space to the brain' they have been able to make connections between how neurons behave in space, and in many non-spatial tasks. They have provided new computational explanations for how these processes might support behaviour, and for the link between space and memory. Going forward, it is an exciting time to be in the field. These contributions have been made by many researchers across the globe, and have relied on a genuine link between theory and experiment. We believe that this cross-disciplinary collaboration is poised to make big strides in understanding how our brains make sense of the structure of our experience, and use it to construct new flexible behaviours.

## References

1. Scoville, W. B. & Milner, B. Loss of recent memory after bilateral hippocampal lesions. *J.Neurol. Neurosurg.Psychiat.* **20**, 11–21 (1957).

2. Cohen, N. J. & Squire, L. R. Preserved learning and retention of pattern-analyzing skill in amnesia: Dissociation of knowing how and knowing that. *Science* **210**, 207–210, 10.1126/science.7414331 (1980).

3. O'Keefe, J. & Nadel, L. *The Hippocampus as a Cognitive Map* (Oxford University Press, 1978).

4. Hafting, T., Fyhn, M., Molden, S., Moser, M.-b. B. & Moser, E. I. Microstructure of a spatial map in the entorhinal cortex. *Nature* **436**, 801–806, 10.1038/nature03721 (2005).

5. Hassabis, D., Kumaran, D., Vann, S. D. & Maguire, E. A. Patients with hippocampal amnesia cannot imagine new experiences. *Proc. Natl. Acad. Sci.* **104**, 1726–1731, 10.1073/pnas.0610561104 (2007).

6. Tolman, E. C. Cognitive maps in rats and men. *Psychol. Rev.* **55**, 189–208, 10.1037/h0061626 (1948).

7. Turner, C. H. The homing of ants: An experimental study of ant behavior. *J. Comp. Neurol. Psychol.* **17**, 10.1002/cne.920170502 (1907).

8. Zanforlin, M. & Poli, G. The Burrowing Rat: A New Technique to Study Place Learning and Orientation. *Acti. et Mem. dell'Academia Patovina di Scienze, Lett. ed Arti* **82**, 653–670 (1970).

9. Behrens, T. E. J. *et al.* What Is a Cognitive Map? Organizing Knowledge for Flexible Behavior. *Neuron* **100**, 490–509, 10.1016/j.neuron.2018.10.002 (2018).

10. Humboldt, W. *On Language: On the Diversity of Human Language Construction and its Influence on the Mental Development of the Human Species* (Cambridge University Press, 0).

11. Tenenbaum, J. B., Kemp, C., Griffiths, T. L. & Goodman, N. D. How to grow a mind: Statistics, structure, and abstraction. *Science* **331**, 1279–1285, 10.1126/science.1192788 (2011).

12. Bartlett, F. C. & Burt, C. Remembering: a Study in Experimental and Social Psychology. *Br. J. Educ. Psychol.* **3**, 187–192, 10.1111/j.2044-8279.1933.tb02913.x (1932).

13. Harlow, H. F. The formation of learning sets. *Psychol. Rev.* **56**, 51–65, 10.1037/h0062474 (1949).

14. Moser, E. I., Moser, M.-B. & McNaughton, B. L. Spatial representation in the hippocampal formation: a history. *Nat. Neurosci.* **20**, 1448–1464, 10.1038/nn.4653 (2017).

15. Aronov, D., Nevers, R. & Tank, D. W. Mapping of a non-spatial dimension by the hippocampalâĂŞentorhinal circuit. *Nature* **543**, 719–722, 10.1038/nature21692 (2017).

16. Knudsen, E. B. & Wallis, J. D. Closed-Loop Theta Stimulation in the Orbitofrontal Cortex Prevents Reward-Based Learning. *Neuron* **106**, 537–547, 10.1016/j.neuron.2020.02.003 (2020).

17. Nieh, E. H. *et al.* Geometry of abstract learned knowledge in the hippocampus. *Nature* 10.1038/s41586-021-03652-7 (2021).

18. Constantinescu, A. O. *et al.* Organizing conceptual knowledge in humans with a gridlike code. *Science* **352**, 1464–1468, 10.1126/science.aaf0941 (2016).

19. Bongioanni, A. *et al.* Activation and disruption of a neural mechanism for novel choice in monkeys. *Nature* **591**, 270–274, 10.1038/s41586-020-03115-5 (2021).

20. Bao, X. *et al.* Grid-like Neural Representations Support Olfactory Navigation of a Two-Dimensional Odor Space. *Neuron* **102**, 1066–1075, 10.1016/j.neuron.2019.03.034 (2019).

21. Park, S. A., Miller, D. S., Nili, H., Ranganath, C. & Boorman, E. D. Map Making: Constructing, Combining, and Inferring on Abstract Cognitive Maps. *Neuron* **107**, 1226–1238, 10.1016/j.neuron.2020.06.030 (2020).

22. Radulescu, A., Shin, Y. S. & Niv, Y. Human Representation Learning. *Annu. Rev. Neurosci.* **44**, 253–273, 10.1146/annurev-neuro-092920-120559 (2021).

23. Sanders, H., Wilson, M. A. & Gershman, S. J. Hippocampal remapping as hidden state inference. *eLife* **9**, e51140, 10.7554/eLife.51140 (2020).

24. Stoianov, I., Maisto, D. & Pezzulo, G. The hippocampal formation as a hierarchical generative model supporting generative replay and continual learning. *bioRxiv preprint* https://doi.org/10.1101/2020.01.16.908889 (2020).

25. O'Keefe, J. & Dostrovsky, J. The hippocampus as a spatial map. Preliminary evidence from unit activity in the freely-moving rat. *Brain Res.* **34**, 171–175, 10.1016/0006-8993(71)90358-1 (1971).

26. Bellmund, J. L. S., Gärdenfors, P., Moser, E. I. & Doeller, C. F. Navigating cognition: Spatial codes for human thinking. *Science* **362**, eaat6766, 10.1126/science.aat6766 (2018).

27. Whittington, J. C. R. *et al.* The Tolman-Eichenbaum Machine: Unifying Space and Relational Memory through Generalization in the Hippocampal Formation. *Cell* **183**, 1249–1263, 10.1016/j.cell.2020.10.024 (2020).

28. Carpenter, F., Manson, D., Jeffery, K., Burgess, N. & Barry, C. Grid cells form a global representation of connected environments. *Curr. Biol.* **25**, 1176–1182, 10.1016/j.cub.2015.02.037 (2015).

29. Wood, E. R., Dudchenko, P. A., Robitsek, R. J. & Eichenbaum, H. Hippocampal neurons encode information about different types of memory episodes occurring in the same location. *Neuron* **27**, 623–633, 10.1016/S0896-6273(00)00071-4 (2000).

30. Niv, Y. Learning task-state representations. *Nat. Neurosci.* **22**, 1544–1553, 10.1038/s41593-019-0470-8 (2019).

31. Bahdanau, D. *et al.* Systematic Generalization: What Is Required and Can It Be Learned? *arXiv preprint* 1–16 (2018).

32. Sutton, R. S. & Barto, A. G. Reinforcement learning: an introduction. *UCL,Computer Sci. Dep. Reinf. Learn. Lect.* 1054, 10.1109/TNN.1998.712192 (2017).

33. Bellman, R. Markovian decision processes (1957).

34. Gershman, S. J. & Niv, Y. Learning latent structure: Carving nature at its joints. *Curr. Opin. Neurobiol.* **20**, 251–256, 10.1016/j.conb.2010.02.008 (2010).

35. Wilson, R. C., Takahashi, Y. K., Schoenbaum, G. & Niv, Y. Orbitofrontal cortex as a cognitive map of task space. *Neuron* **81**, 267–278, 10.1016/j.neuron.2013.11.005 (2014).

36. Watkins, C. J. & Dayan, P. Technical Note: Q-Learning. *Mach. Learn.* **8**, 279–292, 10.1023/A:1022676722315 (1992).

37. Tolman, E. C., Ritchie, B. F. & Kalish, D. Studies in spatial learning. I. Orientation and the short-cut. *J. Exp. Psychol.* **36**, 13–24, 10.1037/h0053944 (1946).

38. Silver, D. *et al.* A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. *Science* **362**, 1140–1144, 10.1126/science.aar6404 (2018).

39. Botvinick, M. & Toussaint, M. Planning as inference. *Trends Cogn. Sci.* **16**, 485–488, 10.1016/j.tics.2012.08.006 (2012).

40. Bush, D., Barry, C., Manson, D. & Burgess, N. Using Grid Cells for Navigation. *Neuron* **87**, 507–520, 10.1016/j.neuron.2015.07.006 (2015).

41. Stemmler, M., Mathis, A. & Herz, A. V. M. Connecting multiple spatial scales to decode the population activity of grid cells. *Sci. Adv.* **1**, e1500816, 10.1126/science.1500816 (2015).

42. Foster, D. J., Morris, R. G. & Dayan, P. A model of hippocampally dependent navigation, using the temporal difference learning rule. *Hippocampus* **10**, 1–16, 10.1002/(SICI)1098-1063(2000)10:1<1::AID-HIPO1>3.0.CO;2-1 (2000).

43. Gustafson, N. J. & Daw, N. D. Grid Cells, Place Cells, and Geodesic Generalization for Spatial Reinforcement Learning. *PLoS Comput. Biol.* **7**, e1002235, 10.1371/journal.pcbi.1002235 (2011).

44. Dayan, P. Improving Generalization for Temporal Difference Learning: The Successor Representation. *Neural Comput.* **5**, 613–624, 10.1162/neco.1993.5.4.613 (1993).

45. Stachenfeld, K. L. K. L. K. L. K. L., Botvinick, M. M. & Gershman, S. J. The hippocampus as a predictive map. *Nat. Neurosci.* **20**, 1643–1653, 10.1038/nn.4650 (2017).

46. Dordek, Y., Soudry, D., Meir, R. & Derdikman, D. Extracting grid cell characteristics from place cell inputs using non-negative principal component analysis. *eLife* **5**, 1–36, 10.7554/eLife.10094 (2016).

47. Mehta, M. R., Quirk, M. C. & Wilson, M. A. Experience-dependent asymmetric shape of hippocampal receptive fields. *Neuron* **25**, 707–715, 10.1016/S0896-6273(00)81072-7 (2000).

48. Derdikman, D. *et al.* Fragmentation of grid cell maps in a multicompartment environment. *Nat. Neurosci.* **12**, 1325–1332, 10.1038/nn.2396 (2009).

49. Krupic, J., Burgess, N. & O'Keefe, J. Neural Representations of Location Composed of Spatially Periodic Bands. *Science* **337**, 853–857, 10.1126/science.1222403 (2012).

50. Garvert, M. M., Dolan, R. J. & Behrens, T. E. E. E. A map of abstract relational knowledge in the human hippocam-palâĂŞentorhinal cortex. *eLife* **6**, 1–20, 10.7554/eLife.17086 (2017).

51. Schapiro, A. C., Turk-browne, N. B., Botvinick, M. M., Norman, K. A. & Schapiro, A. C. Complementary learning systems within the hippocampus: a neural network modelling approach to reconciling episodic memory with statistical learning. *Philos. Transactions Royal Soc. B: Biol. Sci.* **372**, 20160049, 10.1098/rstb.2016.0049 (2017).

52. Momennejad, I. *et al.* The successor representation in human reinforcement learning. *Nat. Hum. Behav.* **1**, 680–692, 10.1038/s41562-017-0180-8 (2017).

53. Momennejad, I. Learning Structures: Predictive Representations, Replay, and Generalization. *Curr. Opin. Behav. Sci.* **32**, 155–166, 10.1016/j.cobeha.2020.02.017 (2020).

54. Piray, P. & Daw, N. D. Unpredictability vs. volatility and the control of learning. *bioRxiv preprint* 2020.10.05.327007 (2020).

55. Todorov, E. Linearly-solvable Markov decision problems. *Adv. Neural Inf. Process. Syst.* 1369–1376, 10.7551/mitpress/7503.003.0176 (2007).

56. Mark, S., Moran, R., Parr, T., Kennerley, S. W. & Behrens, T. E. Transferring structural knowledge across cognitive maps in humans and models. *Nat. Commun.* **11**, 1–12, 10.1038/s41467-020-18254-6 (2020).

57. Dusek, J. A. & Eichenbaum, H. The hippocampus and memory for orderly stimulus relations. *Proc. Natl. Acad. Sci.* **94**, 7109–7114, 10.1073/pnas.94.13.7109 (1997).

58. Frank, L. M., Brown, E. N. & Wilson, M. Trajectory encoding in the hippocampus and entorhinal cortex. *Neuron* **27**, 169–178, 10.1016/S0896-6273(00)00018-0 (2000).

59. Komorowski, R. W., Manns, J. R. & Eichenbaum, H. Robust Conjunctive Item-Place Coding by Hippocampal Neurons Parallels Learning What Happens Where. *J. Neurosci.* **29**, 9918–9929, 10.1523/JNEUROSCI.1378-09.2009 (2009).

60. Grieves, R. M., Wood, E. R. & Dudchenko, P. A. Place cells on a maze encode routes rather than destinations. *eLife* **5**, 1–24, 10.7554/eLife.15986 (2016).

61. Sun, C., Yang, W., Martin, J. & Tonegawa, S. Hippocampal neurons represent events as transferable units of experience. *Nat. Neurosci.* **23**, 651–663, 10.1038/s41593-020-0614-x (2020).

62. George, D. *et al.* Clone-structured graph representations enable flexible learning and vicarious evaluation of cognitive maps. *Nat. Commun.* **12**, 2392, 10.1038/s41467-021-22559-5 (2021).

63. Cormack, G. V. & Horspool, R. N. Data compression using dynamic markov modelling. *Comput. J.* **30**, 541–550, 10.1093/comjnl/30.6.541 (1987).

64. Dempster, A. P., Laird, N. M. & Rubin, D. B. Maximum Likelihood from Incomplete Data Via the EM Algorithm . *J. Royal Stat. Soc. Ser. B (Methodological)* **39**, 1–22, 10.1111/j.2517-6161.1977.tb01600.x (1977).

65. Taube, J., Muller, R. & Ranck, J. Head-direction cells recorded from the postsubiculum in freely moving rats. I. Description and quantitative analysis. *The J. Neurosci.* **10**, 420–435, 10.1523/JNEUROSCI.10-02-00420.1990 (1990).

66. Darwin, C. Origin of certain instincts. *Nature* **7**, 417–418, 10.1038/007417a0 (1873).

67. Mittelstaedt, M. L. & Mittelstaedt, H. Homing by path integration in a mammal. *Naturwissenschaften* **67**, 566–567, 10.1007/BF00450672 (1980).

68. Etienne, A. S. & Jeffery, K. J. Path integration in mammals. *Hippocampus* **14**, 180–192, 10.1002/hipo.10173 (2004).

69. Loomis, J. M. *et al.* Nonvisual navigation by blind and sighted: Assessment of path integration ability. *J. Exp. Psychol. Gen.* **122**, 73–91, 10.1037/0096-3445.122.1.73 (1993).

70. Maaswinkel, H., Jarrard, L. E. & Whishaw, I. Q. Hippocampectomized rats are impaired in homing by path integration. *Hippocampus* **9**, 553–561, 10.1002/(SICI)1098-1063(1999)9:5<553::AID-HIPO9>3.0.CO;2-G (1999).

71. Sreenivasan, S. & Fiete, I. Grid cells generate an analog error-correcting code for singularly precise neural computation. *Nat. Neurosci.* **14**, 1330–1337, 10.1038/nn.2901 (2011).

72. Mathis, A., Herz, A. V. & Stemmler, M. Optimal population codes for space: Grid cells outperform place cells. *Neural Comput.* **24**, 2280–2317, 10.1162/NECO{_}a{_}00319 (2012).

73. Chen, G., Lu, Y., King, J. A., Cacucci, F. & Burgess, N. Differential influences of environment and self-motion on place and grid cell firing. *Nat. Commun.* **10**, 1–11, 10.1038/s41467-019-08550-1 (2019).

74. Zhang, K. Representation of spatial orientation by the intrinsic dynamics of the head-direction cell ensemble: a theory. *The J. Neurosci.* **16**, 2112–2126, 10.1523/JNEUROSCI.16-06-02112.1996 (1996).

75. Skaggs, W. E., Knierim, J. J., Kudrimoti, H. S. & McNaughton, B. L. A model of the neural basis of the rat's sense of direction. *Adv. neural information processing systems* **7**, 173–180 (1995).

76. Samsonovich, A. & McNaughton, B. L. Path integration and cognitive mapping in a continuous attractor neural network model. *J. Neurosci.* **17**, 5900–5920, 10.1523/jneurosci.17-15-05900.1997 (1997).

77. Tsodyks, M. Attractor neural network models of spatial maps in hippocampus. *Hippocampus* **9**, 481–489, 10.1002/(SICI)1098-1063(1999)9:4<481::AID-HIPO14>3.0.CO;2-S (1999).

78. Burak, Y. & Fiete, I. R. Accurate path integration in continuous attractor network models of grid cells. *PLoS Comput. Biol.* **5**, e1000291, 10.1371/journal.pcbi.1000291 (2009).

79. Ben-Yishai, R., Bar-Or, R. L. & Sompolinsky, H. Theory of orientation tuning in visual cortex. *Proc. Natl. Acad. Sci.* **92**, 3844–3848, 10.1073/pnas.92.9.3844 (1995).

80. Kim, S. S., Rouault, H., Druckmann, S. & Jayaraman, V. Ring attractor dynamics in the Drosophila central brain. *Science* **356**, 849–853, 10.1126/science.aal4835 (2017).

81. Yoon, K. *et al.* Specific evidence of low-dimensional continuous attractor dynamics in grid cells. *Nat. Neurosci.* **16**, 1077–1084, 10.1038/nn.3450 (2013).

82. Gardner, R. J. *et al.* Toroidal topology of population activity in grid cells. *bioRxiv preprint* 2021.02.25.432776 (2021).

83. O'Keefe, J. & Recce, M. L. Phase relationship between hippocampal place units and the EEG theta rhythm. *Hippocampus* **3**, 317–330, 10.1002/hipo.450030307 (1993).

84. Burgess, N., Barry, C. & O'Keefe, J. An oscillatory interference model of grid cell firing. *Hippocampus* **17**, 801–812, 10.1002/hipo.20327 (2009).

85. Banino, A. *et al.* Vector-based navigation using grid-like representations in artificial agents. *Nature* **557**, 429–433, 10.1038/s41586-018-0102-6 (2018).

86. Cueva, C. J. & Wei, X.-X. Emergence of grid-like representations by training recurrent neural networks to perform spatial localization. *Int. Conf. on Learn. Represent.* **0**, 1–19 (2018).

87. Sorscher, B., Mel, G. C., Ganguli, S. & Ocko, S. A. A unified theory for the origin of grid cells through the lens of pattern formation. *Adv. Neural Inf. Process. Syst. 32* **32**, 10003–10013 (2019).

88. Gentner, D. Structure-mapping: A theoretical framework for analogy. *Cogn. Sci.* **7**, 155–170, 10.1016/S0364-0213(83)80009-3 (1983).

89. Anderson, M. I. & Jeffery, K. J. Heterogeneous modulation of place cell firing by changes in context. *J. Neurosci.* **23**, 8827–8835 (2003).

90. Bostock, E., Muller, R. U. & Kubie, J. L. Experience-dependent modifications of hippocampal place cell firing. *Hippocampus* **1**, 193–205, 10.1002/hipo.450010207 (1991).

91. Muller, R. U. & Kubie, J. L. The effects of changes in the environment on the spatial firing of hippocampal complex-spike cells. *J. Neurosci.* **7**, 1951–1968 (1987).

92. Fyhn, M., Hafting, T., Treves, A., Moser, M. B. & Moser, E. I. Hippocampal remapping and grid realignment in entorhinal cortex. *Nature* **446**, 190–194, 10.1038/nature05601 (2007).

93. Carey, S. & Bartlett, E. Acquiring a single new word. *Pap. Reports on Child Lang. Dev.* **15**, 17–29 (1978).

94. Manns, J. R. & Eichenbaum, H. Evolution of declarative memory. *Hippocampus* **16**, 795–808, 10.1002/hipo.20205 (2006).

95. Kemp, C. & Tenenbaum, J. B. The discovery of structural form. *Proc. Natl. Acad. Sci.* **105**, 10687–10692, 10.1073/pnas.0802631105 (2008).

96. Lake, B. M., Salakhutdinov, R. & Tenenbaum, J. B. Human-level concept learning through probabilistic program induction. *Science* **350**, 1332–1338, 10.1126/science.aab3050 (2015).

97. Battleday, R. M. & Griffiths, T. L. Analogy as nonparametric bayesian inference over relational systems. *arXiv preprint* (2020).

98. Ellis, K. *et al.* DreamCoder: Growing generalizable, interpretable knowledge with wake-sleep Bayesian program learning. *arXiv preprint* 1–22 (2020).

99. Høydal, Ã. A., Skytøen, E. R., Andersson, S. O., Moser, M.-B. & Moser, E. I. Object-vector coding in the medial entorhinal cortex. *Nature* **568**, 400–404, 10.1038/s41586-019-1077-7 (2019).

100. Hartley, T., Burgess, N., Lever, C., Cacucci, F. & O'Keefe, J. Modeling place fields in terms of the cortical inputs to the hippocampus. *Hippocampus* **10**, 369–379, 10.1002/1098-1063(2000)10:4<369::AID-HIPO3>3.0.CO;2-0 (2000).

101. Becker, S. & Burgess, N. Modelling spatial recall, mental imagery and neglect. In *Advances in Neural Information Processing Systems 13*, 96–102 (2001).

102. Barry, C. *et al.* The boundary vector cell model of place cell firing and spatial memory. *Rev. Neurosci.* **17**, 71–97, 10.1515/REVNEURO.2006.17.1-2.71 (2006).

103. Solstad, T., Boccara, C. N., Kropff, E., Moser, M.-B. & Moser, E. I. Representation of Geometric Borders in the Entorhinal Cortex. *Science* **322**, 1865–1868, 10.1126/science.1166466 (2008).

104. Lever, C., Burton, S., Jeewajee, A., O'Keefe, J. & Burgess, N. Boundary vector cells in the subiculum of the hippocampal formation. *J. Neurosci.* **29**, 9771–9777, 10.1523/JNEUROSCI.1319-09.2009 (2009).

105. Gauthier, J. L. & Tank, D. W. A Dedicated Population for Reward Coding in the Hippocampus. *Neuron* **99**, 179–193, 10.1016/j.neuron.2018.06.008 (2018).

106. Sarel, A., Finkelstein, A., Las, L. & Ulanovsky, N. Vectorial representation of spatial goals in the hippocampus of bats. *Science* **355**, 176–180, 10.1126/science.aak9589 (2017).

107. Uria, B. *et al.* The spatial memory pipeline: A model of egocentric to allocentric understanding in mammalian brains. *bioRxiv preprint* 1–52, 10.1101/2020.11.11.378141 (2020).

108. Pritzel, A. *et al.* Neural Episodic Control. *arXiv preprint* 10.1038/nature20101 (2017).

109. Hebb, D. O. *The Organization of Behavior; A Neuropsychological Theory* (Wiley and Sons, New York, 1949).

110. Hopfield, J. J. Neural networks and physical systems with emergent collective computational abilities (associative memory/parallel processing/categorization/content-addressable memory/fail-soft devices). *Biophysics* **79**, 2554–2558 (1982).

111. McKenzie, S. *et al.* Hippocampal representation of related and opposing memories develop within distinct, hierarchically organized neural schemas. *Neuron* **83**, 202–215, 10.1016/j.neuron.2014.05.019 (2014).

112. Bunsey, M. & Eichenbaum, H. Conservation of hippocampal memory function in rats and humans. *Nature* **379**, 255–257, 10.1038/379255a0 (1996).

113. Sanders, H., Wilson, M., Klukas, M., Sharma, S. & Fiete, I. Efficient Inference in Structured Spaces. *Cell* **183**, 1147–1148, 10.1016/j.cell.2020.11.008 (2020).

114. Grieves, R. M. & Jeffery, K. J. The representation of space in the brain, 10.1016/j.beproc.2016.12.012 (2017).

115. Eichenbaum, H. Time cells in the hippocampus: A new dimension for mapping memories. *Nat. Rev. Neurosci.* **15**, 732–744, 10.1038/nrn3827 (2014).

116. Teyler, T. J. & Rudy, J. W. The hippocampal indexing theory and episodic memory: Updating the index. *Hippocampus* **17**, 1158–1169, 10.1002/hipo.20350 (2007).

117. Marr, D. Simple memory: a theory for archicortex. *Philos. transactions Royal Soc. London. Ser. B, Biol. sciences* **262**, 23–81, 10.1098/rstb.1971.0078 (1971).

118. McClelland, J. L., McNaughton, B. L. & O'Reilly, R. C. Why there are complementary learning systems in the hippocampus and neocortex: Insights from the successes and failures of connectionist models of learning and memory. *Psychol. Rev.* **102**, 419–457, 10.1037/0033-295X.102.3.419 (1995).

119. Evans, T. & Burgess, N. Coordinated hippocampal-entorhinal replay as structural inference. *Adv. Neural Inf. Process. Syst. 32* **32**, 1731–1743 (2019).

120. Boccara, C. N., Nardin, M., Stella, F., OâĂŹNeill, J. & Csicsvari, J. The entorhinal cognitive map is attracted to goals. *Science* **363**, 1443–1447, 10.1126/science.aav4837 (2019).

121. Friston, K. The free-energy principle: a rough guide to the brain? *Trends Cogn. Sci.* **13**, 293–301, 10.1016/j.tics.2009.04.005 (2009).

122. Levine, S. Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review. *arXiv preprint* arXiv:1805.00909v3 (2018).

123. McNamee, D. C., Stachenfeld, K. L., Botvinick, M. M. & Gershman, S. J. Flexible modulation of sequence generation in the entorhinalâĂŞhippocampal system. *Nat. Neurosci.* 10.1038/s41593-021-00831-7 (2021).

124. Sims, D. W. *et al.* Hierarchical random walks in trace fossils and the origin of optimal search behavior. *Proc. Natl. Acad. Sci. United States Am.* **111**, 11073–11078, 10.1073/pnas.1405966111 (2014).

125. Pfeiffer, B. E. & Foster, D. J. Autoassociative dynamics in the generation of sequences of hippocampal place cells. *Science* **349**, 180–183, 10.1126/science.aaa9633 (2015).

126. Baram, A. B., Muller, T. H., Whittington, J. C. R. & Behrens, T. E. Intuitive planning: global navigation through cognitive maps based on grid-like codes. *bioRxiv preprint* **0**, 421461, 10.1101/421461 (2018).

127. Yu, C., Behrens, T. E. & Burgess, N. Prediction with directed transitions: complex eigenstructure, grid cells and phase coding. *arXiv preprint* 1–21 (2020).

128. Burak, Y. & Fiete, I. Do we understand the emergent dynamics of grid cell activity? *J. Neurosci.* **26**, 9352–9354, https://doi.org/10.1523/JNEUROSCI.2857-06.2006 (2006).

129. Wang, J. X. *et al.* Prefrontal cortex as a meta-reinforcement learning system. *Nat. Neurosci.* **21**, 860–868, 10.1038/s41593-018-0147-8 (2018).

130. Duan, Y. *et al.* RL$^2$: Fast Reinforcement Learning via Slow Reinforcement Learning. *arXiv preprint* 1–14, 10.1051/0004-6361/201527329 (2016).

131. Chen, L. *et al.* Decision Transformer: Reinforcement Learning via Sequence Modeling. *arXiv preprint* 1–19 (2021).

132. Janner, M., Li, Q. & Levine, S. Reinforcement Learning as One Big Sequence Modeling Problem. *arXiv preprint* (2021).

133. Foster, D. J. & Wilson, M. A. Reverse replay of behavioural sequences in hippocampal place cells during the awake state. *Nature* **440**, 680–683, 10.1038/nature04587 (2006).

134. Deshmukh, S. S. & Knierim, J. J. Influence of local objects on hippocampal representations: Landmark vectors and memory. *Hippocampus* **23**, 253–67, 10.1002/hipo.22101 (2013).

135. Mattar, M. G. & Daw, N. D. Prioritized memory access explains planning and hippocampal replay. *Nat. Neurosci.* **21**, 1609–1617, 10.1038/s41593-018-0232-z (2018).

136. Ólafsdóttir, H. F., Carpenter, F. & Barry, C. Coordinated grid and place cell replay during rest. *Nat. Neurosci.* **19**, 792–794, 10.1038/nn.4291 (2016).

137. Kaefer, K., Nardin, M., Blahna, K. & Csicsvari, J. Replay of Behavioral Sequences in the Medial Prefrontal Cortex during Rule Switching. *Neuron* **106**, 154–165, 10.1016/j.neuron.2020.01.015 (2020).

138. Butler, W. N., Hardcastle, K. & Giocomo, L. M. Remembered reward locations restructure entorhinal spatial maps. *Science* **363**, 1447–1452, 10.1126/science.aav5297 (2019).

139. Ziv, Y. *et al.* Long-term dynamics of CA1 hippocampal place codes. *Nat. Neurosci.* **16**, 264–266, 10.1038/nn.3329 (2013).

140. Driscoll, L. N., Pettit, N. L., Minderer, M., Chettih, S. N. & Harvey, C. D. Dynamic Reorganization of Neuronal Activity Patterns in Parietal Cortex. *Cell* **170**, 986–999, 10.1016/j.cell.2017.07.021 (2017).

141. Rule, M. E., O'Leary, T. & Harvey, C. D. Causes and consequences of representational drift. *Curr. Opin. Neurobiol.* **58**, 141–147, 10.1016/j.conb.2019.08.005 (2019).

142. Rubin, A., Geva, N., Sheintuch, L. & Ziv, Y. Hippocampal ensemble dynamics timestamp events in long-term memory. *eLife* **4**, 1–16, 10.7554/eLife.12247 (2015).

143. Pastalkova, E., Itskov, V., Amarasingham, A. & Buzsaki, G. Internally Generated Cell Assembly Sequences in the Rat Hippocampus. *Science* **321**, 1322–1327, 10.1126/science.1159775 (2008).

144. MacDonald, C. J., Lepage, K. Q., Eden, U. T. & Eichenbaum, H. Hippocampal "time cells" bridge the gap in memory for discontiguous events. *Neuron* **71**, 737–749, 10.1016/j.neuron.2011.07.012 (2011).

145. Chen, T., Kornblith, S., Norouzi, M. & Hinton, G. A Simple Framework for Contrastive Learning of Visual Representations. *arXiv preprint* (2020).

146. Salz, D. M. *et al.* Time cells in hippocampal area CA3. *J. Neurosci.* **36**, 7476–7484, 10.1523/JNEUROSCI.0087-16.2016 (2016).

147. Higgins, I. *et al.* β-VAE: Learning basic visual concepts with a constrained variational framework. *Int. Conf. on Learn. Represent.* **0** (2017).

148. Zhou, J., Jia, C., Montesinos-cartagena, M. & Gardner, M. P. H. Evolving schema representations in orbitofrontal ensembles during learning. *Nature* 10.1038/s41586-020-03061-2 (2020).

149. Zhou, J. *et al.* Complementary Task Structure Representations in Hippocampus and Orbitofrontal Cortex during an Odor Sequence Task. *Curr. Biol.* **29**, 3402–3409, 10.1016/j.cub.2019.08.040 (2019).

150. Miller, E. K. & Cohen, J. D. An integrate theory of PFC function. *Annu. Rev. Neurosci.* **24**, 167–202 (2001).

151. Bernardi, S. *et al.* The Geometry of Abstraction in the Hippocampus and Prefrontal Cortex. *Cell* **183**, 954–967, 10.1016/j.cell.2020.09.031 (2020).

152. Morton, N. W., Schlichting, M. L. & Preston, A. R. Representations of common event structure in medial temporal lobe and frontoparietal cortex support efficient inference. *Proc. Natl. Acad. Sci. United States Am.* **117**, 29338–29345, 10.1073/pnas.1912338117 (2020).

153. Samborska, V., Butler, J. L., Walton, M. E., Behrens, T. E. & Akam, T. Complementary task representations in hippocampus and prefrontal cortex for generalising the structure of problems. *bioRxiv preprint* **4**, 2021.03.05.433967 (2021).

154. Schuck, N. W., Cai, M. B., Wilson, R. C. & Niv, Y. Human Orbitofrontal Cortex Represents a Cognitive Map of State Space. *Neuron* **91**, 10.1016/j.neuron.2016.08.019 (2016).

155. Yu, J. Y., Liu, D. F., Loback, A., Grossrubatscher, I. & Frank, L. M. Specific hippocampal representations are linked to generalized cortical representations in memory. *Nat. Commun.* **9**, 1–11, 10.1038/s41467-018-04498-w (2018).

156. Ito, H. T., Zhang, S. J., Witter, M. P., Moser, E. I. & Moser, M. B. A prefrontal-thalamo-hippocampal circuit for goal-directed spatial navigation. *Nature* **522**, 50–55, 10.1038/nature14396 (2015).

157. Hawkins, J., Lewis, M., Klukas, M., Purdy, S. & Ahmad, S. A Framework for Intelligence and Cortical Function Based on Grid Cells in the Neocortex. *Front. Neural Circuits* **12**, 1–15, 10.3389/fncir.2018.00121 (2019).

158. Lewis, M. Hippocampal Spatial Mapping As Fast Graph Learning. *arXiv preprint* (2021).

159. Hochreiter, S. & Schmidhuber, J. Long Short-term Memory. *Neural Comput.* **9**, 17351780 (1997).

160. Vaswani, A. *et al.* Attention is all you need. *Adv. Neural Inf. Process. Syst.* **2017-Decem**, 5999–6009 (2017).

161. Brown, T. B. *et al.* Language models are few-shot learners. *arXiv preprint* (2020).

162. Dosovitskiy, A. *et al.* An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *arXiv preprint* 1–21 (2020).

163. Dehaene, S., Meyniel, F., Wacongne, C., Wang, L. & Pallier, C. The Neural Representation of Sequences: From Transition Probabilities to Algebraic Patterns and Linguistic Trees. *Neuron* **88**, 2–19, 10.1016/j.neuron.2015.09.019 (2015).

164. Christiansen, M. H. & Chater, N. Toward a connectionist model of recursion in human linguistic performance. *Cogn. Sci.* **23**, 157–205, 10.1207/s15516709cog2302{_}2 (1999).

165. Amalric, M. & Dehaene, S. Origins of the brain networks for advanced mathematics in expert mathematicians. *Proc. Natl. Acad. Sci. United States Am.* **113**, 4909–4917, 10.1073/pnas.1603205113 (2016).

166. Nieder, A. Supramodal numerosity selectivity of neurons in primate prefrontal and posterior parietal cortices. *Proc. Natl. Acad. Sci. United States Am.* **109**, 11860–11865, 10.1073/pnas.1204580109 (2012).

167. Higgins, I. *et al.* Towards a Definition of Disentangled Representations. *arXiv preprint* 1–29 (2018).

168. Killian, N. J. & Buffalo, E. A. Grid cells map the visual world. *Nat. Neurosci.* **21**, 161–162, 10.1038/s41593-017-0062-4 (2018).

169. Nau, M., Navarro Schröder, T., Bellmund, J. L. S. & Doeller, C. F. Hexadirectional coding of visual space in human entorhinal cortex. *Nat. Neurosci.* **21**, 10.1038/s41593-017-0050-8 (2018).

170. Julian, J. B., Keinath, A. T., Frazzetta, G. & Epstein, R. A. Human entorhinal cortex represents visual space using a boundary-anchored grid. *Nat. Neurosci.* **21**, 191–194, 10.1038/s41593-017-0049-1 (2018).

171. Schwartenbeck, P. *et al.* Generative replay for compositional visual understanding in the prefrontal-hippocampal circuit. *bioRxiv preprint* 2021.06.06.447249 (2021).

## Acknowledgements

## Author Contributions

JCRW and TEJB conceptualised manuscript. JCRW and DM performed simulations. JCRW and DM drafted manuscript. JCRW and TEJB edited manuscript with input from all other authors.

## Data availability

No data was generated in this perspective.

## Code availability

Python and Tensorflow code will made available on publication.