



## COGNITIVE NEUROSCIENCE

# A multilevel account of hippocampal function in spatial and concept learning: Bridging models of behavior and neural assemblies

Robert M. Mok<sup>1\*</sup> and Bradley C. Love<sup>2,3\*</sup>

A complete neuroscience requires multilevel theories that address phenomena ranging from higher-level cognitive behaviors to activities within a cell. We propose an extension to the level of mechanism approach where a computational model of cognition sits in between behavior and brain: It explains the higher-level behavior and can be decomposed into lower-level component mechanisms to provide a richer understanding of the system than any level alone. Toward this end, we decomposed a cognitive model into neuron-like units using a neural flocking approach that parallels recurrent hippocampal activity. Neural flocking coordinates units that collectively form higher-level mental constructs. The decomposed model suggested how brain-scale neural populations coordinate to form assemblies encoding concept and spatial representations and why so many neurons are needed for robust performance at the cognitive level. This multilevel explanation provides a way to understand how cognition and symbol-like representations are supported by coordinated neural populations (assemblies) formed through learning.

## INTRODUCTION

Neuroscience is a multilevel enterprise. Its target of explanation ranges from behavioral to molecular phenomena. Satisfying and complete explanations of the mind and brain will necessarily be multilevel (1–3). In multilevel componential (or constitutive) explanations, each component at a higher level can be decomposed into its own lower-level mechanism (1). For example, the circulatory system's capacity to deliver oxygen and energy to the body can be decomposed into lower-level mechanisms including the heart's blood pumping and kidney's blood filtering mechanism, which together supports the function of the circulatory system. These mechanisms themselves can be decomposed into their components, such as the muscle contractions of the heart or filtering units of the kidney, which, in turn, can be further decomposed as desired (1).

Marr's (4) well-known three level organization is one multilevel proposal but seems inappropriate for neuroscience as it relegates all of neuroscience to one level, namely, the implementational level, whereas neuroscience itself is a multilevel endeavor. We extend Craver (1) by proposing a mechanistic multilevel approach in which the top level is behavior and the first level of mechanism below behavior is an algorithmic model that captures behavior and whose components can be related to brain measures. The components of this mechanism can be further decomposed into their own mechanisms to address finer-grain scientific questions (e.g., neural populations; Fig. 1).

Under this account, the highest-level mechanism could be a cognitive model that captures behavior. The components of this cognitive model could be related to neural measures and further decomposed. What is a component at a higher level is a mechanism

at a lower level that can itself be decomposed into components to account for additional findings and make new predictions. For a mechanism to be truly multilevel, lower-level decompositions must fully reproduce phenomena at higher levels, up to and including behavior. This successful decomposition into constituent mechanisms is what provides explanatory power in multilevel accounts. The power of this approach is that mechanisms at different levels are not unrelated efforts aiming to explain different types of data. Instead, multilevel explanations can be integrated and offer a more complete understanding of the domain in which one can "zoom in or out" to the level of granularity desired for the current question.

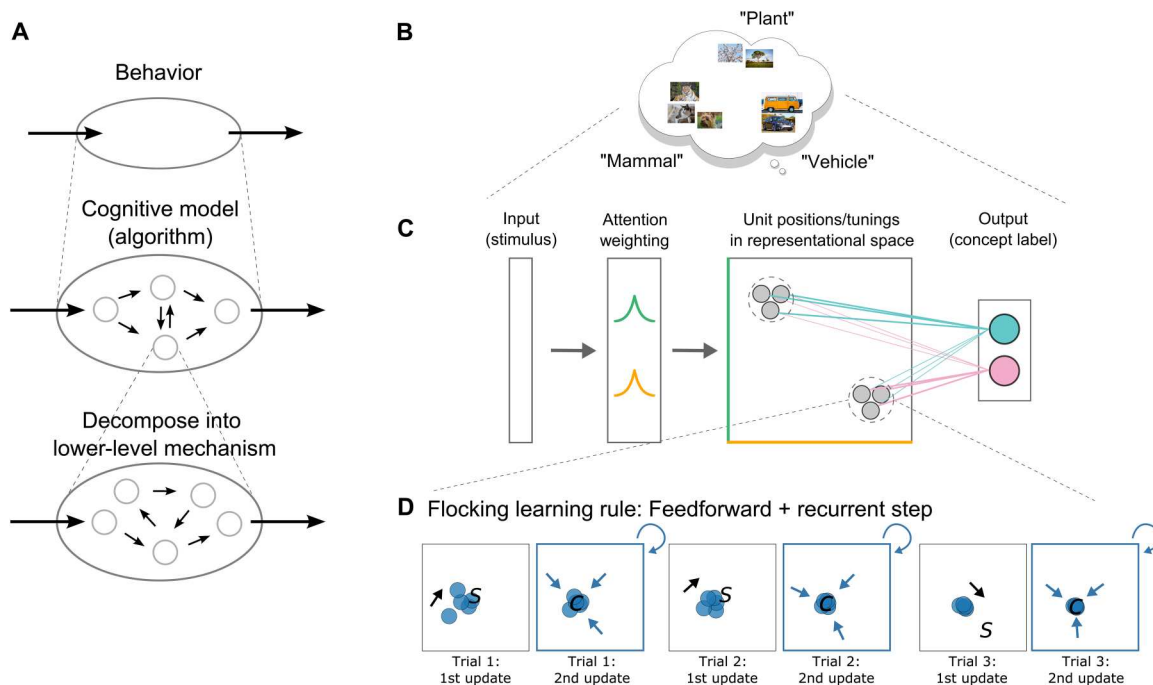
Constructing theories and testing their predictions at multiple levels provide a more comprehensive account of the system of interest. Although neuroscience is guilty of a bias toward lower-level explanations (5, 6), higher-level mechanisms are crucial in multilevel explanation because they offer explanatory concepts not available at lower levels (3). For example, the heart's contractions make little sense without considering the function of the circulatory system, and the hippocampal synaptic weights and activity patterns make little sense without notions of memory and learning. In neuroscience, it is common to construct a specific theory or model to fit the specific data at hand, which unfortunately leads to a disconnected patchwork of theories for individual phenomena. Multilevel theories can weave this patchwork together into a coherent and complete account in which each level is situated within mechanisms that lie above and below. As one descends levels, additional phenomena can be addressed, whereas, as one ascends, the function of the mechanism within the overall system becomes clearer.

Neuroscience has very few multilevel theories of this variety that can bridge between cognitive constructs and neuronal activity—multilevel explanations from behavior to neurons. For example, how does the brain implement a symbol? Specifically, how do brain systems coordinate neural populations to form symbol-like representations, such as highly selective neural assemblies (7–9) that encode concepts (10) or specific spatial locations (11)? One

<sup>1</sup>MRC Cognition and Brain Sciences Unit, University of Cambridge, 15 Chaucer Road, Cambridge CB2 7EF, UK. <sup>2</sup>UCL Department of Experimental Psychology, 26 Bedford Way, London WC1H 0AP, UK. <sup>3</sup>The Alan Turing Institute, London, United Kingdom.

\*Corresponding author. Email: rob.mok@mrc-cbu.cam.ac.uk (R.M.M.); b.love@ucl.ac.uk (B.C.L.)

Copyright © 2023  
The Authors, some  
rights reserved;  
exclusive licensee  
American Association  
for the Advancement  
of Science. No claim to  
original U.S. Government  
Works. Distributed  
under a Creative  
Commons Attribution  
License 4.0 (CC BY).



**Fig. 1. Multilevel explanation of concept learning in the brain: decomposition of a cognitive model into neural flocks.** (A) Levels of mechanisms for neuroscience. Behavior is the phenomenon of interest, explained by a task-performing algorithm (cognitive model), which can be decomposed into lower-level mechanisms. (B) Our behavior of interest is concept learning and categorization. (C) Behavior is explained by the cognitive model. After the stimulus is encoded, attention is applied, and neuron-like units activate according to their similarity to the input. These activations are transmitted through learned association weights to generate an output (e.g., category decision). Dotted circles are abstract clusters. (D) Decomposition of the clusters into neuron-like units and the flocking learning rule. Clusters from (C) are decomposed into neuron-like units [gray circles in (C) represent units, and dashed circles highlight units in the same flock or virtual cluster]. Left to right:  $k$  winners (blue) move toward the stimulus ("S"), followed by a recurrent update, where units move toward their centroid ("C").  $k$  neuron-like units become similarly tuned over time, forming a neural flock.

suggestion is that the hippocampal formation represents physical space (12–14) and abstract spaces for encoding concepts and abstract variables (15–17), constructing cognitive maps (18) for mental navigation in these spaces. However, it is unclear how populations of similarly tuned neurons in the hippocampus acquire their highly selective tuning properties to concepts or spatial locations. While tantalizing, identifying a cell tuned to particular concept, such as Jennifer Aniston (10), does not specify the supporting mechanisms leaving such findings as curiosities that invite explanation.

Attempts have been made to offer multilevel theories in neuroscience, but critical explanatory gaps remain. In our own work, we have developed a cognitive model, SUSTAIN, of how people learn categories from examples. SUSTAIN addresses a number of behavioral findings (19, 20), and aspects of the model have been related to the brain (20–23). The hippocampus was related to SUSTAIN's clustering mechanism, which bundles together relevant information in memory during learning, and this account was verified by a number of brain imaging studies (20–22). The goal-directed attentional mechanism in SUSTAIN was linked to the ventral medial prefrontal cortex, and this account was verified by brain imaging (23) and animal lesion studies (24). These same mechanisms provide a good account of place and grid cell activity in spatial tasks (17).

Although successful in addressing both behavior and accompanying brain activity, this line of work, similar to almost all work in model-based neuroscience, is limited to (i) proposing

correspondences between model components and brain regions and (ii) evaluating these correspondences in terms of model-brain correlates. However, how do we move beyond mere neural correlates to a lower-level mechanistic explanation that unpacks the higher-level theory? Cognitive models, such as SUSTAIN, come with abstract constructs such as clusters, and it is left entirely open how they could be decomposed into the neural populations that give rise to behavior. It is insufficient to state that each cognitive construct (e.g., a cluster) is instantiated by a number of neurons, just as it is unsatisfying to state that Jennifer Aniston is somehow represented by multiple neurons, the spreadsheet on a computer relies on a number of transistors, and so forth. How do neurons coordinate to give rise to the higher-level construct? This is the key question that needs to be addressed to move beyond mere neural correlates toward a mechanistic understanding of how the brain implements cognition.

We aim to address this explanatory gap by decomposing aspects of a cognitive model, SUSTAIN, into a mechanism consisting of neuron-like units. Critically, the aggregate action of these neuron-like units give rise to virtual structures akin to the higher-level cognitive constructs in SUSTAIN (i.e., clusters) while retaining SUSTAIN's account of behavior (Fig. 1, C and D). By taking a key component of a cognitive model that addresses behavior and decomposing it to the level of a neuron, we offer an account of how concept cells and spatially tuned cell assemblies can arise in the brain.

One of the main challenges is how to bridge from abstract constructs such as clusters to neurons, while retaining the high-level behavior of the model. How do single neurons track a concept? How does the brain coordinate the activity of many neurons during learning, as there are thousands of neurons, but only a few clusters are required to represent a concept at the cognitive level? That is, how does a select population of neurons learn and become tuned to similar features in the world such as in concept cells and place cells rather than independently develop their own tuning? To implement a higher-level construct, such as a cluster or symbol, neurons must somehow solve this coordination problem.

Inspired by algorithms that capture flocking behavior in birds and fish (25), we propose that the hippocampus may exhibit neural flocking, where coordinated activity—virtual clusters or flocks—arises from local rules (Fig. 1D). This coordination could be achieved by recurrent processing in the hippocampus [e.g., (26)], which we formalize in learning rules in which neuron-like units that are highly active in response to a stimulus (e.g., a photo of Jennifer Aniston) both adjust their tuning toward that stimulus and to each other. The latter learning update is the key to flocking as members of the neural flock coordinate by moving closer to each other. That simple addition to standard learning rules is sufficient to solve the coordination problem and give rise to virtual clusters and hippocampal place or concept cell (10) assemblies (8).

Gazing at the neuron-like units forming the model, one will not see clusters just as one will not see clusters nor symbols by peering into the gray goo of the brain. Nevertheless, the coordinated activity of these neuron-like units can be described as supporting these higher-level constructs that behave in aggregate like the higher-level cognitive model, SUSTAIN. In the model specification and simulations that follow, we aim to close the aforementioned explanatory gap and make the case for multilevel explanation in neuroscience. In addition, by decomposing the higher-level model, we can consider how the brain benefits in terms of fault and noise tolerance by implementing clusters in a large neuronal flock. For instance, the mental representation of a concept is preserved when one neuron, or even a set of neurons in the neural assembly, dies, as well as when there is synaptic turnover over relatively short time scales (typical in the hippocampus) (27). Finally, we consider how the model can be extended and further decomposed to account for differences in processing across anterior-posterior hippocampus axis.

## RESULTS

### Multineuron clustering model

Our multineuron clustering model, SUSTAIN-d, is a decomposition of the SUSTAIN model of category learning. Whereas prototype models always form one unit in memory for each category and exemplar models store one unit for each episode, SUSTAIN moves between these two extremes depending on the nature of the learning problem. SUSTAIN assumes that the environment is regular and clusters similar experiences together in memory until there is an unexpected prediction error, such as encountering a bat and wrongly predicting that it is a bird based on its similarity to an existing bird cluster. When such an error occurs, a new cluster is recruited. Notable episodes (e.g., encountering a bat for the first time) can transition to concepts over time (e.g., other similar bats are stored in the same cluster). The hippocampus is critical in supporting this

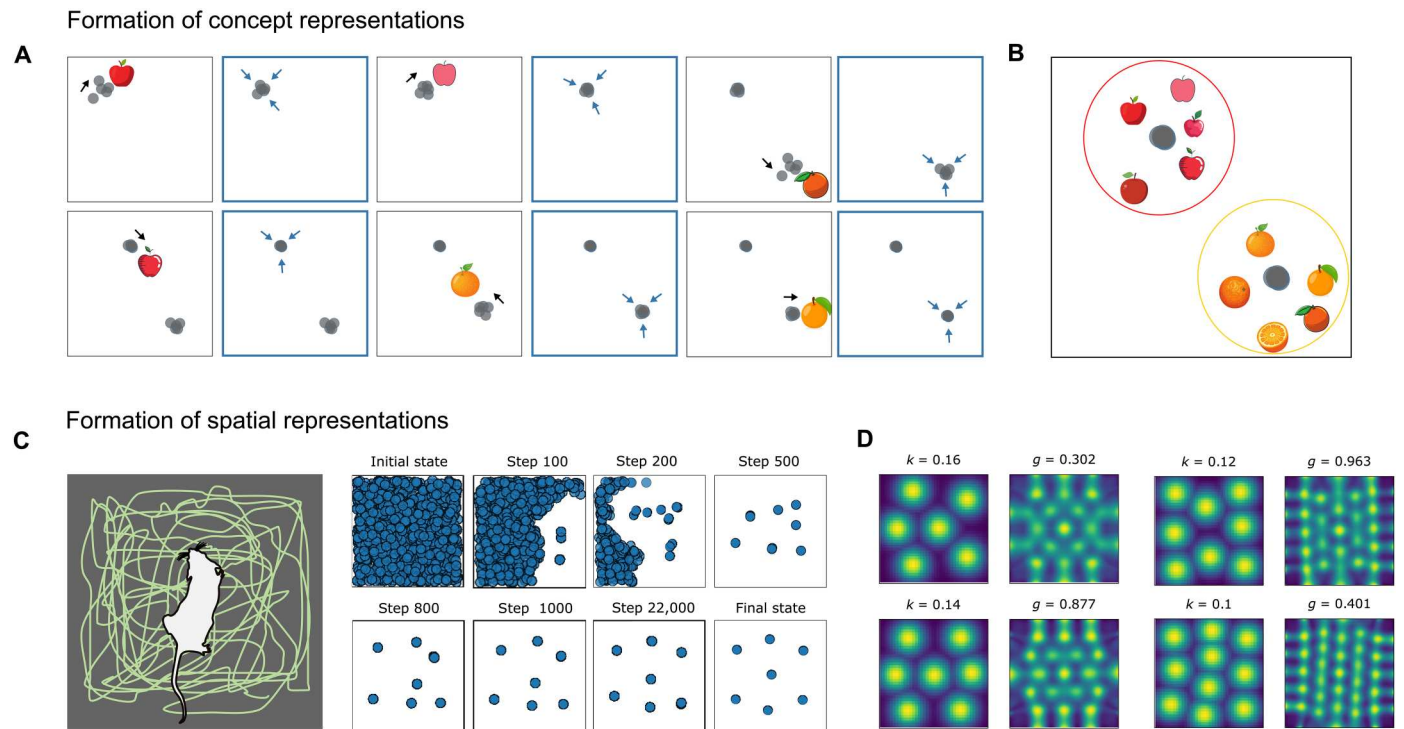
form of learning (28). SUSTAIN has other elements that are not the focus of this contribution, such as an attentional mechanism that determines which aspects of stimuli are most relevant to categorization.

Here, we decompose SUSTAIN into more neuron-like units while retaining its overall function. We will refer to this model as SUSTAIN-d for SUSTAIN decomposed (see Materials and Methods for formal model description). Both models use a local learning rule, the Kohonen learning rule, but SUSTAIN-d has a second update rule that implements flocking through recurrence. Both models have a recruitment mechanism, but rather than recruit cognitive units like clusters, SUSTAIN-d, like the hippocampus, has an existing pool of neuron-like computing that enter each task with some tuning (i.e., preferentially activated for particular stimuli). Unlike SUSTAIN, which will recruit a handful of clusters for a learning problem, SUSTAIN-d can consist of an arbitrarily large number of computing units (see below for brain-scale simulations where the number of computing units is equal to the number of neurons in the hippocampus). Finally, both models have attention weights with a local update rule and have connection weights from clusters or units to outputs that are updated through task-related error.

Despite these notable differences, SUSTAIN-d's thousands of neuron-like computing units show the same aggregate behavior as SUSTAIN. This is accomplished by solving the aforementioned coordination problem by what we refer to as neural flocking (Fig. 1D). The key to neural flocking is that units that are highly activated by a stimulus both adjust their tunings toward the stimulus and each other. This double update leads to virtual clusters forming that can consist of thousands of neuron-like computing units. In general, the number of clusters SUSTAIN recruits will match the number of neural flocks that arise in SUSTAIN-d, which leads to the models providing equivalent behavioral accounts. Whereas SUSTAIN associates clusters with a category or response, SUSTAIN-d's individual neuron-like units form connection weights (Fig. 1C). In summary, SUSTAIN-d is formulated absent of cognitive constructs like clusters, but, nevertheless, its neuron-like units in aggregate manifest the same principles and can account for the same behaviors, which provides an account of how cognitive constructs can be implemented in the brain.

### Formation of concept and spatial representations by neural flocking

SUSTAIN-d's neural flocking mechanism can explain how concept cells (Fig. 2, A and B) and spatially tuned place and grid cells arise (Fig. 2, C and D). Whereas our previous work (17, 19) relied on higher-level mental constructs (i.e., clusters) to account for such phenomena, here, we show how a neural population can coordinate to virtually form such structures via the flocking learning rule. In the spatial domain, we simulated an agent (e.g., rodent) exploring a square environment (free foraging), which leads to virtual clusters akin to place cells distributed in a grid-like pattern (Fig. 2C), which, in turn, leads to cells that monitor these units' activity displaying a grid-like firing pattern (Fig. 2D). An alternative model with no recurrence (i.e., no flocking rule) shows no self-organization of spatial cells (fig. S5A). In the conceptual domain, where the representation space is not as uniformly sampled, virtual clusters are clumpier (Fig. 2B), and monitoring units will show no grid response. These results in the conceptual and spatial domain hold across a wide



**Fig. 2. Formation of concept and spatial representations by neural flocking.** (A) The model learns distinct representations for apples and oranges.  $k$  winners (i.e., most activated units) adjust their receptive fields toward the current stimulus, followed by a recurrent update toward their centroid. (B) The second update is sufficient to solve the coordination problem allowing SUSTAIN-d to form neural flocks or virtual clusters (which, in this example, represent the concepts apple and orange). (C) Spatial representation formation. Left: An agent (e.g., a rodent) forages in an environment. Right: Development of spatial representations. SUSTAIN-d's neuron-like units are initially uniformly tuned to locations. At each time step, the  $k$  winners move toward the stimulus (e.g., sensory information at the current location) and each other (i.e., neural flocking). This learning dynamic creates flocks or virtual clusters of units with similar spatial tuning, akin to place-cell assemblies. These flocks tile the environment. (D) Examples of grid cell-like activity patterns and corresponding spatial autocorrelograms after learning. See fig. S1A for more examples and fig. S1B for distributions of grid scores).

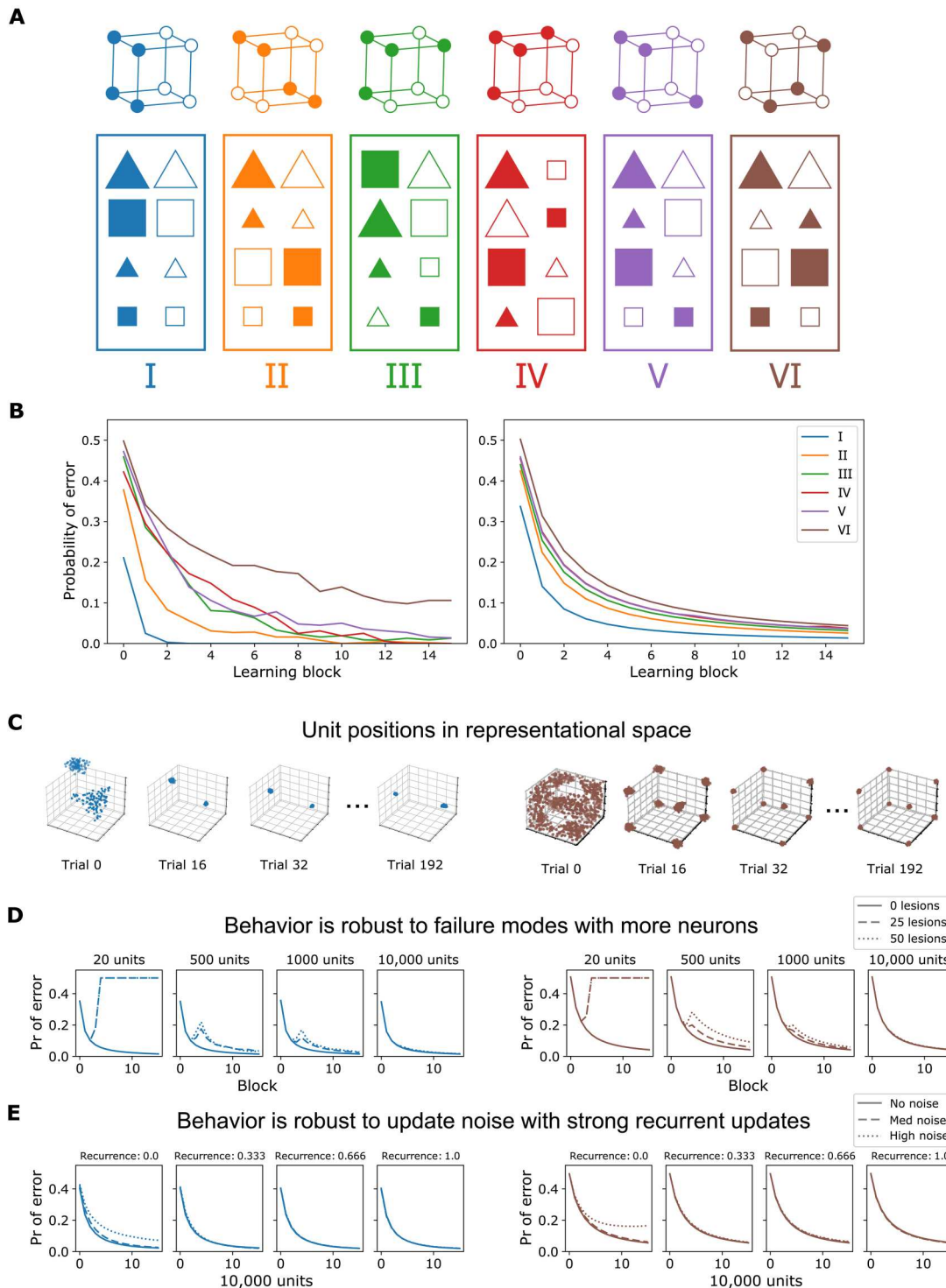
range of parameters. The flocks that arise from the interaction of numerous neuron-like units are a higher-level representation in that the number of underlying neuron-like units involved can vary over orders of magnitude with the aggregate behavior at the "flock level" remaining constant. This robustness allows SUSTAIN-d to decompose cognitive constructs to neuron-like units at same scale as the brain as demonstrated by simulations where the number of units is equal to the number of neurons in the corresponding brain regions.

### Neural population-based model retains high-level cognitive model properties and captures concept learning behavior

One major challenge for our multilevel proposal is to account for complex behaviors that hitherto were the sole province of cognitive models. Can SUSTAIN-d with its neuron-like units account for the same behaviors that SUSTAIN does by relying on cognitive constructs? We evaluate whether SUSTAIN-d can account for human learning performance on the six classic learning problems of Shepard *et al.* (29) (Fig. 3). To provide a true multilevel theory, we aim for SUSTAIN-d's solution in terms of virtual clusters arising from neural flocking to parallel SUSTAIN's clustering solutions, which provide a good correspondence to hippocampal activity patterns (22).

SUSTAIN-d was trained in a manner analogous to human participants, learning through trial and error in these supervised learning tasks. On each trial, SUSTAIN-d updated its neuron-like units' positions in representational space, attention weights, and connections weights from its neuron-like to category responses (see Materials and Methods for details). All the learning updates were error-driven and local, as opposed to propagating errors over multiple network layers as in deep learning models. Whereas SUSTAIN forms a new cluster in response to an unexpected error (e.g., learning that a bat is not a bird), SUSTAIN-d recruits the  $k$ -nearest unconnected neuron-like units to the current stimulus, which is in accord with the intuition that the brain repurposes existing resources during learning. The neural flocking learning rule leads to these recruited units forming a virtual cluster.

SUSTAIN-d captured the difficulty ordering of the human learning curves (Fig. 3B, right), and its solutions paralleled those of SUSTAIN in terms of the modal number of clusters recruited (two, four, six, six, six, and eight flocks for each of the six learning problems) and attentional allocation to features. Notably, SUSTAIN-d's results scale to a large number of neuron-like units, producing the same output and learning curves from few (e.g., 50) to many neurons [ $3.2 \times 10^6$  hippocampal principal cells (30, 31) as used here]. Thus, SUSTAIN-d provides a multilevel account (1) of hippocampally mediated learning that ranges from behavior to



**Fig. 3. SUSTAIN-d's brain-scale population of neuron-like units collectively displays the same behavior as the high-level cognitive model that it decomposes, while making additional predictions about robustness in neural computation. (A)** Six concept learning structures (29). Bottom: In each box, stimuli in the left and right columns are in different categories. Top: Cubes represent each stimulus in binary stimulus feature space (color, shape, and size) for each structure. **(B)** Learning curves from human behavior (65) (left) and model fits (right). Probability of error is plotted as a function of learning block for each structure. **(C)** Neuron-like units form neural flocks or virtual clusters (e.g., type I in blue and type VI in brown; see fig. S2 for all types) that parallel clusters in the higher-level cognitive model. The number of units are subsampled from the whole population for better visualization. **(D)** The more neuron-like units, the more robust the model is when confronted by failure modes (e.g., cell death, noise, and synaptic transmission failure). **(E)** The stronger the recurrence during learning, the better the noise tolerance. See fig. S3 for more examples.

Downloaded from https://www.science.org at University College London on July 31, 2023

neuron-like units. SUSTAIN-d is able to display similar aggregate behavior over a wide range of neuron-like units because its learning updates and operations can be scaled to reflect the number of units involved (see Materials and Methods).

Similar to the human brain, SUSTAIN-d is resistant to minor insults and faults. Each neural assembly or flock can consist of many neuron-like units (Fig. 3C and fig. S2), not all of which are needed for the remaining units to function as a virtual cluster. SUSTAIN can be viewed as SUSTAIN-d when the number of highly activated units in response to the stimulus is 1 (i.e.,  $k = 1$ ). As  $k$  or total number of units increase, lesioning a subset of SUSTAIN-d's units has negligible effects on aggregate behavior (Fig. 3D and fig. S3). This robustness through redundancy is consistent with operation of the brain where multiple neurons and synapses with similar properties offset the effects of damage and natural turnover of dendritic spines (27, 32–34).

Having multiple units combined with SUSTAIN-d's recurrent update can also counteract noise (Fig. 3E). In these simulations, we added noise to SUSTAIN-d's neuron-like units, which will lead to units from other assemblies (or neural flocks) becoming highly active, which can lead to incorrect decisions and disrupt learning. SUSTAIN's recurrent update (Fig. 1, C and D) ameliorates these effects of noise by pulling individual units' responses toward the mean of flock (Fig. 3E and fig. S4). The same self-organizing learning rules that enable SUSTAIN-d's neuron-like units to behave in aggregate like a cognitive model also make it more robust to noise and damage.

Furthermore, recurrence is particularly important when coordinating large populations of noisy neurons, as the resulting combined magnitude of noise is markedly stronger. We find that the recurrence strength markedly reduces the effect of noise for 10,000-unit models compared to 20-unit models (fig. S4). Finally, alternative models such as a prototype model (one unit per concept label) are not capable of capturing learning behavior and the disruptive effects of lesions and noise (fig. S5B).

### Further decomposing to capture differential function in anterior and posterior hippocampus

In a multilevel mechanistic account, model components can be further decomposed to capture findings that require more fine-grain mechanisms. SUSTAIN-d decomposed SUSTAIN's clusters into neuron-like units. Here, we further decompose SUSTAIN-d's neuron-like units into two pools to capture functional differences between anterior and posterior hippocampus.

Anterior place fields tend to be broader and lower granularity than posterior hippocampal fields (35), and this appears to be a general principle at the population level (36, 37). For category learning studies, one prediction is that learning problems with a simpler structure that promotes broad generalization would be more anterior, whereas learning problems that have a complex irregular structure would be better suited to posterior hippocampus. This pattern holds across studies (21, 22) (Fig. 4A). Here, we simulate Shepard's six learning problems, which order from what should be most anterior (type I) to most posterior (type VI). Type I can be solved by focusing and generalizing on the basis of one stimulus feature, whereas type VI requires memorizing the details of each item.

Although the anterior-posterior distinction may best be viewed as a continuous axis, we simplified to create two banks of neuron-like units for SUSTAIN-d, one corresponding to anterior

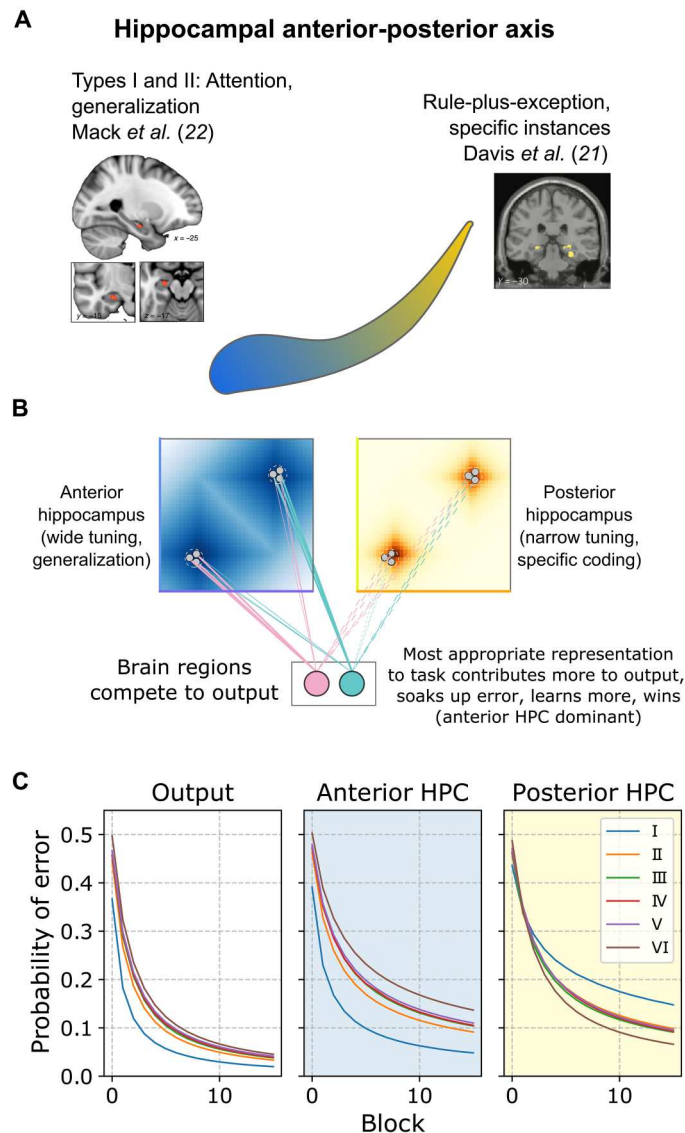
hippocampus and one to posterior hippocampus. These two banks of neuron-like units only differed in how broad their tuning was with anterior fields being broader than posterior fields. The responses from both pools of neuron-like units were pooled to determine the category response (Fig. 4B). The two pools of units formed neural flocks independently and were only linked in that they both aimed to correctly predict the category label. Thus, the pool that was more suited to a learning problem will take over by developing larger connection weights to the output units indicating the response. In effect, both pools of units compete to guide the overall response (see Materials and Methods). Thus, the anterior pool of units should guide the overall system response when a problem is very simple (e.g., type I), whereas the posterior pool should take over when a problem is highly irregular and complex (e.g., type VI).

As predicted, the anterior pool guided the overall response for simple problems like the type I problem, whereas the posterior pool did so for complex problems, like the type VI problem (Fig. 4C). The posterior pool learns the type VI problem faster than the type I problem because SUSTAIN-d forms eight virtual cluster (neural flocks) for the type VI problem, which is ideally suited to the narrow fields in the posterior pool that favor memorization over broad generalization, exactly the opposite functional properties of the anterior pool. These simulations suggest that the anterior-posterior axis in the hippocampus may provide complementary learning mechanisms.

### DISCUSSION

Neuroscience is inherently a multilevel enterprise that requires multilevel theories to move from a patchwork understanding of the mind and brain toward more integrative and encompassing theories. One solution, which we advance, is to adopt a levels of mechanism approach (1) in which the components of a higher-level model are unpacked into their own mechanisms. Specifically, we extended the levels of mechanism approach to computational models, effectively combining it with levels of analysis (4) and building an account that bridges across levels from behavior to algorithm to multilevel mechanisms. This is particularly important for cognitive neuroscience, where many mechanisms are best described via an abstract algorithmic account, and the challenge is to specify how it could arise from lower-level mechanisms. Notably, this is not constrained to a single implementational level, as one can selectively decompose arbitrarily lower levels for the aspects of the model that are of scientific interest depending on the question at hand, from macroneural systems to neural assemblies and neurotransmitters, down to ion channels. It has been assumed that levels can be linked, but there has been no concrete effort to specify how—we provided a solution on how to bridge across relevant levels in neuroscience. Here, we demonstrate this through a multilevel explanation of concept learning in which a cognitive model at the top level that accounts for behavior is decomposed into populations of neuron-like units that form assemblies.

This breakthrough was largely achieved by a novel learning rule that promoted neural flocking in which neuron-like units with related receptive field properties became more similar by coordinating using recurrence. Similar to a flock of birds, these units functioned as a collective, providing an account of how symbols and other cognitive constructs, such as clusters, can arise from



**Fig. 4. Further decomposing SUSTAIN-d to capture known functional differences along the anterior-posterior axis of the hippocampus.** (A) Illustration of the anterior-posterior (blue-yellow) gradient in human hippocampus (HPC). Place fields in the anterior hippocampus are broader, and posterior hippocampus place fields are more narrow. Likewise, anterior hippocampus is strongly activated by concept learning structures that follow broad, general rules (left), and posterior hippocampus is more strongly engaged by irregular rule-plus-exception structures where specific instances are important. (B) SUSTAIN-d is further decomposed into a bank of units with broader tuning to model anterior hippocampus (blue) and a narrowly tuned bank of units to model posterior hippocampus (yellow). Both banks contribute to the output and compete to exert control over the category decision. (C) Model output (left) combines the anterior (middle) and posterior (right) neuron-like units' output. The anterior units dominate for simple category structures, whereas the posterior units dominate for irregular structures.

neuron-like computing elements following biologically consistent operations.

The recurrent update in the learning rule was inspired by biological recurrence in the hippocampal formation (within the hippocampus and big-loop recurrence in the medial temporal lobe) (26), where multiple passes allow for deeper processing. Symbol-like neural representations naturally form through our implementation of recurrence, suggesting that functional neural assemblies can form through a flocking mechanism like in birds. With recent advances in large-scale neural recordings over time, future neurophysiological or imaging studies in animals that record from hippocampus across time and over learning could search for such a neural mechanism.

With a neural population and a recurrent mechanism, the model also naturally captures the brain's tendency to encode the same information in many neurons (i.e., redundancy) (32–34), which makes the system more tolerant to neuronal damage, natural synaptic turnover, and noise. Notably, there could be more than one recurrent update, as there are multiple recurrent loops in the brain. Future work could introduce more recurrent steps or different forms of recurrence such as constraining it by known anatomical pathways.

Note that the model assumes a sparse representation where most neurons are dormant but active neurons are highly active for a small group of stimuli, which parallels sparse coding in the hippocampus. Furthermore, the sparse code of a fixed proportion of highly activated winners in a neural population is consistent with recent hippocampal neurophysiological evidence showing that the overall activity of place cells within an environment is constant (38). This work shows how localist coding models (i.e., with sparse, highly tuned neurons) can be implemented in a neural population (39–42). One way for related approaches to showing mixed selectivity is for different stimulus dimensions to be relevant in different contexts. For example, a neuron may respond to wool socks in a clothing context and to tart flavors in a food context with the relevant dimensions or features in each context being nonoverlapping.

Our model makes several experimental predictions. One prediction is that neural flocking in the hippocampus will lead to receptive fields harmonizing over learning. Specifically, cells initially activated in response to novel stimuli will form a neural assembly, and their tunings will become more harmonized over stimulus repetitions (forming a virtual cluster). According to the theory, it will be (largely) the same group of neurons that will keep responding [k-winners-take-all (k-WTA)], consistent with sparse coding in the hippocampus. This could be tested by recording many neurons in rodents or nonhuman primates during trial-by-trial learning. It would be particularly interesting to record across multiple regions to test for recurrent processing across the hippocampal-cortical loop and how this corresponds to our implementation of recurrence. Our flocking mechanism suggests that initially activated cells to novel stimuli should flock together and become highly tuned neural assemblies such as concept and place cell assemblies.

Another prediction relates to the stability of the hippocampus (43) including representational changes due to synaptic turnover (27) and neuron death (44). Despite these changes, hippocampal representations appear to be relatively consistent over time [e.g., concepts cells (10) and place cells (45, 46)], and the hippocampus appears to be important for precise long-term memories (47). In our model, a group of neuron-like units are recruited when the

task demands it. As demonstrated in the lesion simulations, if the model only has a few neuron-like units, then there is a marked detriment to performance when units are removed (e.g., from neuron death), but with a larger population of units, the model is more resilient to such insults. When the magnitude of losses is sufficiently large, the model will make incorrect decisions, which leads to recruitment of units to maintain performance. This account is consistent with the relative stability in the hippocampus for representing long-term concepts and locations, while, at the same time, having instability in the biological substrate [see (48, 49) for a potential mechanism for place or concept cell recruitment]. Hence, the model predicts that new neurons will be recruited if a sufficient number of neurons lose their synaptic efficacy, die, or tuning simply drifts. This could be tested by tracking neural populations over time with calcium imaging or new neurophysiological techniques that can record the same neurons over long periods of time (50).

The hippocampus exhibits mixed selectivity where neurons respond to different features across contexts [e.g., (51)]. In the tasks we considered, the feature space was low dimensional such that only one set of attention weights was required. However, the model could be extended such that neuron-like units would display a unique peak response for each context. Each context would correspond to a different subspace with its own set of attention weights. Future work will explore these approaches to capturing mixed selectivity in hippocampal cells and consider how this tuning supports behavior.

One benefit of multilevel theories is that model components can be decomposed into their own mechanisms as desired. In effect, one can selectively zoom in to consider the aspects of the mechanism of interest at a finer grain. In the last set of simulations, we further decomposed SUSTAIN-d's neuron-like units into two banks of units corresponding to anterior and posterior hippocampus. Further decomposing SUSTAIN-d allows us to account for finer-grain phenomena and make new predictions.

We found that the bank of units best suited to the task dominated learning. For learning problems with a simple structure, the broader receptive fields of the anterior bank dominated. When the learning problem had an irregular structure that required memorization, the posterior bank dominated. By varying the broadness of the receptive fields, we introduced a generally applicable framework in which modules compete with one another to guide behavior but are ultimately complementary in terms of accomplishing the goals of the overall system. This cross-regional competitive framework captures the common finding across cognitive neuroscience where brain regions that have the more appropriate or useful representations for the task at hand are more strongly activated and contribute more to behavior.

In the future, we plan to extend the method to include connections across modules (e.g., excitatory/inhibitory) based on known anatomy and functional properties and to model more interacting brain regions. The neuron-like units themselves could be further decomposed and elaborated to behave more like biological neurons. For our present purposes, we did not require spiking neurons with complex dynamics. However, just as SUSTAIN was decomposed into SUSTAIN-d, so too can SUSTAIN-d be decomposed as desired, all the while retaining the overarching principles and behavior of the higher-level models.

In our work, we used the Kohonen learning rule that is simple, based on local computation, and has powerful information-processing capabilities. However, the rule is still abstract, and future work can decompose it into a lower-level mechanism with biologically constrained implementational details. Furthermore, there are different ways to consider the sources of error (52–54) and learning mechanisms such as temporal-difference error-driven learning in complementary learning system models of the hippocampus (55–57). Future work can consider different learning rules and their decompositions into lower-level mechanisms to assess how the hippocampus learns from error.

We can extend this approach to other functions of the hippocampus and neocortex such as episodic memory. The role of the hippocampus in concept learning is closely related its role in episodic memory formation, where trial-by-trial learning leads to concept formation “one episode at a time” (28). Future work could construct multilevel accounts for theories of hippocampal function for memory, such as transitive inference (58), statistical learning (59), and consolidation (47), to assess how populations of neurons across multiple brain regions learn under these contexts.

Our account of the hippocampal formation is an alternate view to prior models that focus on navigation and path integration [e.g., (60)] or structure learning (61, 62). In prior work (17), we proposed that hippocampal cells including concept and place cells play a representational role, whereas medial entorhinal cortex (mEC) grid and nongrid spatial cells monitor the activity of hippocampal cells and play an error-monitoring role for cluster or cell recruitment. mEC cells contain information as to whether there exists a hippocampal place or concept cell that encode the current location or stimulus, and if error is high (no field in mEC cells), then the hippocampus can recruit a new cell to represent this new experience. Despite this, it is possible that grid-like representations in our model could be reused for other functions. Our model produces spatial cell-like representations, and so downstream brain areas can use this location information for path integration, but this is a consequence a general hippocampal learning algorithm. More generally, our work aims to address the domain-general properties of the hippocampal formation across concept learning, concept representations, and aspects of spatial representations, and our current model extends this an account that provides an explanation to the level of neural assemblies.

In sum, cognitive neuroscience can benefit from multilevel explanations by exploring and bridging mechanisms across levels. We have many cognitive models that characterize behavior successfully but are in need to be decomposed into a set of mechanistic processes that could be implemented in the brain. In recent years, neuroscience is finally putting more emphasis on behavior (6), but we suggest that for a complete account of the cognitive function of interest, a successful high-level explanation of the behavior (e.g., through a cognitive model) that can be decomposed into the relevant lower-level mechanisms is key.

## MATERIALS AND METHODS

### Overview and motivation of the model

SUSTAIN-d is a decomposition of SUSTAIN (19), a cognitive model of concept learning that has captured behavior in a number of tasks (19, 20) and brain-activity patterns including in the hippocampus and medial temporal lobe structures (17, 21, 22)



[also see (63)]. The formal specification of SUSTAIN is included in the aforementioned papers. Whereas SUSTAIN contains clusters (a cognitive construct) that are recruited in response to unexpected events, SUSTAIN-d decomposes the notion of cluster into neuron-like units that coordinate to form virtual clusters through a flocking learning rule as described below.

**Model learning implementation details**

The model was initialized with a population of neuron-like units. Full-scale simulations used 32,000,000 units to model the hippocampus principal cell population. To determine the best fitting parameters for these simulations (see below), 10,000 units were used to reduce computational costs. Units were placed randomly (uniformly) in the stimulus feature space, where all units are inactive or unconnected to the task context. On the first trial (no output) or when the model makes an error (greater output for the incorrect category),  $k$  (proportion of total; set to 0.00005 or 0.005% for the hippocampus simulation and 0.01 for parameter search) neurons are recruited at the current stimulus' position (note that number of units and  $k$  do not change model behavior; see section on scaling below). Once units are recruited, they are connected and activate in response to stimulus input, and their activations contribute to the category decision. On each trial, the winners' activations contribute to the output decision, and they update their tuning by moving toward the current stimulus (Kohonen learning rule) and then toward their own centroid (recurrent update), and the attention weights and connection (i.e., output) weights are updated through learning.

Specifically, the model takes a stimulus vector as input on each trial; the most strongly activated  $k$  connected neuron-like units are considered winners (k-WTA), and the activation of these units is computed on each trial

$$act_i = \zeta \cdot e^{-\zeta \cdot dist_i} \tag{1}$$

where  $\zeta$  is a positive scaling parameter that controls the steepness of the tuning curve and  $dist_i$  is the attention-weighted distance between neuron  $i$ 's position  $pos_i$  and the stimulus  $x$  in the  $R^n$  representational space they inhabit

$$dist_i = \sum_{j=1}^n [a_j \cdot |pos_{ij} - x_j|^r]^{1/r} \tag{2}$$

where  $r$  is set to 1 for the city-block distance (for separable-dimension stimuli) and the non-negative  $a_j$  attention weights sum to 1. The  $n$  attention weights correspond to each stimulus feature dimension, and their magnitude reflects the importance of each feature. The attention weighting also corresponds to each unit's receptive field, which is centered on its position along each feature dimension (Fig. 1B, attention weighting).  $\zeta$  controls the steepness of the receptive field and will be used to model the different tuning properties of the anterior and posterior hippocampus (see below). Hence, units most strongly tuned to the current stimulus input are activated, and the activity is propagated forward to produce the categorization decision. Only winners have non-zero output that contribute to the category decision

$$out_i = \begin{cases} act_i, & \text{if unit } i \text{ is a winner} \\ 0, & \text{otherwise} \end{cases} \tag{3}$$

and evidence for category decisions propagates from the output

$$evidence_j = \sum_{i=1}^m w_{ij} \cdot out_i \tag{4}$$

where  $w_{ij}$  is the connection weight (Fig. 1B, cyan and pink connections) between unit  $i$  and decision  $j$  and  $m$  is the number of units. Finally, the probability of making decision  $j$  is computed by

$$prob_j = \frac{e^{\phi \cdot evidence_j}}{\sum_{v=1}^z e^{\phi \cdot evidence_v}} \tag{5}$$

where  $z$  is the number of possible decisions and  $\phi$  is a non-negative decision certainty parameter or inverse temperature [see (19, 64) for related formulations]. If there are no connected units (e.g., first trial) or fewer than  $k$  winners, then no units respond or fewer than  $k$  units respond, respectively.

During learning, the  $k$  winners update twice, once toward the stimulus and a second time toward each other, which supports neural flocking. We view unit updating as a continuous process through time relying on recurrent connections, which we simplify here to two simple updates. In the first update, the  $k$  winners update their positions toward the current stimulus' position on each trial according to the Kohonen learning rule

$$\Delta pos_i = \eta_{pos} \cdot (x - pos_i) \tag{6}$$

where  $\eta_{pos}$  is the learning rate,  $x$  is the current stimulus' location, and  $pos_i$  is unit  $i$ 's position in representational space (Fig. 1C, first update). Bold type is reserved for vectors. This is the first learning step, where the initial forward pass of stimulus information occurs and the first clustering update is applied. In the second update, the  $k$  winners perform an additional recurrent step which re-codes the stimulus based on their activity, updating their positions toward the centroid of the  $k$  winners' positions

$$\Delta pos_i = \eta_{group} \cdot (group - pos_i) \tag{7}$$

where  $\eta_{group}$  is the learning rate for the recurrent update and group is the centroid (mean position) of all the winners (Fig. 1C, second update). With this double-update rule, neurons update their tuning profiles to activate more to stimuli that inhabit that part of space in a coordinated fashion, and over the course of learning will typically stabilize into a portion of the space adaptive for the task. As co-activated units cluster together and become more similar to each other in tuning, this group naturally lead to a virtual cluster (i.e., many units that are similar tuned to a concept) similar to a cluster (19) and akin to a hippocampal place or concept cell assembly where assemblies of neurons show similar tuning and coactivate to similar stimuli or environmental features (8-10).

To update attention weights over learning, the model applies gradient ascent on the summed activation of winning units minus the summed activation of the nonwinner units (i.e., connected but not a winner) with learning rate  $\eta_{attn}$ . Note that this learning rule is local (i.e., no backpropagation).

Connection weights are updated using descent on error (based on the output  $prob_j$ ) using cross-entropy loss with one-hot target vector (i.e., one category is correct) on each trial, with learning rate  $\eta_{cweights}$ . To display the output of the model and for fitting to the behavioral data, we plot  $1 - prob_j$  (where  $j$  is the correct category

decision), averaged for each block to match the error learning curves as done in prior work.

As the model activations, outputs, and learning update values will vary depending on the number of units, we scaled the learning rates to retain consistent outputs across models with different numbers of units. This scaling means that changing the total number of units and  $k$  (winners; proportion of total) does not change the learning behavior and output of the model, meaning that the model can scale from a small model with a few neurons to millions of neurons without losing its theoretical essence and cognitive capacities, allowing us to bridge across levels (assuming no noise—see Results for beneficial effects of a larger number of neuron-like units with noise). Furthermore, this allowed us to perform parameter search for fitting human behavior with fewer units and different  $k$ , while reporting the hippocampal-scale simulation with many more units. First, the learning rate of the connection weights was divided by  $k$ , so that weight updates would scale by the  $k$  winners that contribute to the output on each trial. As the attention weights are updated locally by gradient ascent to the winner neurons' activations relative to the loser neurons' activations, we divided the update (the gradient) by the number of active units (i.e., total number of the winner and loser neurons included to compute the gradient).

### Model fitting: Human concept learning behavior

To model the classic Shepard *et al.* (29) results, we fit the learning curve data (minimizing sum of squared errors) from the Nosofsky *et al.* (65) replication of the study of Shepard *et al.* (29). We present a brief overview of this classic study here. Participants learned to categorize eight stimuli that varied on three binary feature dimensions (shape, size, and color) into two categories. The concept structure was one of six possible logical structures from Shepard *et al.* (29) (Fig. 3A). On each trial, participants categorized each stimulus into a category and was provided feedback, learning by trial and error. Participants completed blocks of 16 trials (with two repetitions of each stimulus). Participants continued learning until they made no errors in four sub-blocks of eight trials or if they completed 25 blocks (400 trials). In both studies, they plotted error curves (1 – proportion correct) for the first 16 blocks (16 trials per block), which are the data we will fit (Fig. 3B, left).

Task blocks consisted of 16 stimuli presented in a randomized order. To obtain error curves for each parameter set, a random stimulus sequence for each problem type was generated 25 times, and the error curves were produced by taking the mean across those iterations. To maintain consistency, each iteration was seeded with a specific number, so that the 25 sequences were the same across the different parameters.

Model learning curves were fit to display the human pattern of results. We performed a hierarchical grid search across the parameters. For the standard model (i.e., no separation of brain regions or modules), there were six free parameters:  $\zeta$ ,  $\phi$ , and four learning rates (attention weights  $\eta_{\text{attn}}$ , connection weights  $\eta_{\text{cweights}}$ , Kohonen update  $\eta_{\text{pos}}$ , and recurrent update  $\eta_{\text{group}}$ ). To fit the model that includes a separate anterior and posterior bank of units, separate tuning ( $\zeta$ ) parameters were used for each bank with the constraint that the anterior bank should have a broader tuning than the posterior bank of units (12 free parameters).

### Robustness to failure modes: Noise and lesion experiments

To demonstrate the beneficial effect of having a population of neurons (rather than a single unit or cluster in cognitive models) and the recurrent update, we simulated Shepard's problems with different failure modes during the learning process and how robust the model was to these perturbations. To simulate noise in the learning process, we added noise to the units' positions. For each trial, noise was sampled from a  $n$ -dimensional Gaussian distribution (corresponding to  $n$  features) with zero mean and SD of 0, 0.5, or 1.0 that was added to the update. By adding noise to the unit's position in representational space, this causes potential problems for (i) selecting the appropriate  $k$  neurons as winners, (ii) appropriate updating of the attention weights, and (iii) appropriate updating of the connection weights.

To simulate damage-like events in the neural population, we performed a lesion-like experiment where we randomly removed a subset of the active neurons from the model, simulating typical biological changes such as neuron death or synaptic turnover. For a simple illustration of the beneficial effect of the number of neurons on damage-like events, we set up one "lesion" event at trial 60 where 0, 25, or 50 units were removed and rendered inactive from that point on. The results hold with more lesion events or a larger number of neurons removed.

### Unsupervised learning on spatial tasks

To simulate a rodent foraging in an environment, we placed an agent in a two-dimensional square environment and produced a randomly generated 500,000 sets of steps with the restriction that the agent could not step out of the environment. On each trial, it was able to move left, right, up, or down in steps of 0, 0.025, 0.05, or 0.075. The environment was a square that spanned from 0 to 1 on the horizontal and vertical dimensions.

For unsupervised learning, the model could recruit units like SUSTAIN does by relying on a surprise signal. Here, we further simplify as in (17) and assume that all units in the population are relevant to the current context. Unit positions were updated according to the learning rules specified above. While a rodent's actual environment contains many features, we assumed that these features effectively reduce to a two-dimensional space corresponding to coordinates within the agent's enclosure.

On each trial, the agent moved a step (randomly selected over four directions and four step sizes; one trial), and the model updated the  $k$  winners with the Kohonen learning rule as before, with an annealed learning rate so that the units would eventually settle and stabilize into a particular location

$$\Delta \mathbf{pos}_i = \eta_t \cdot (\mathbf{x} - \mathbf{pos}_i)$$

where  $\eta_t$  is the learning rate at time  $t$ . The learning rate followed an annealing schedule

$$\eta_t = \frac{\eta_0}{1 + \rho \cdot t}$$

where  $\eta_0$  is the initial learning rate and  $\rho$  is the annealing rate set to  $4 \times 10^{-12}$  [see (17)]. The recurrence update learning rate  $\eta_{\text{group}}$  was fixed at 1.0, although smaller values such as 0.8 and 0.6 produce similar results.

To compute grid scores at the end of learning, activation maps were produced by generating a new movement trajectory with

250,000 steps (as above) and computing the unit activations based on their positions at the end of learning [i.e., freezing the positions; see (17)]. For each value of  $k$ , we ran 100 simulations and computed the grid scores at the end of learning. The activation maps were binned in  $40 \times 40$  bins (original  $100 \times 100$ ) and then normalized by the number of visits to each binned location (normalized activation map). Grid scores were calculated on the basis of (66). Briefly, the spatial autocorrelogram of the activation maps were calculated as defined in (67), and gridness was computed using the expanding gridness method, where a circular annulus with a radius of eight bins was placed on the center of the autocorrelogram, with the central peak removed. The annulus was rotated in 30 steps, and the Pearson correlation between the rotated and unrotated version of the spatial autocorrelogram was recorded. The highest correlation value for 30, 90, and 150 rotations was subtracted from the lowest correlation value at 0, 60, and 120 to give an interim grid score. This was repeated expanding the annuli by two bins, up to 20. The final grid score was the highest interim grid score.

### Supplementary Materials

This PDF file includes:

Figs. S1 to S5

[View/request a protocol for this paper from Bio-protocol.](#)

### REFERENCES AND NOTES

- C. F. Craver, *Explaining the Brain* (Oxford Univ. Press, 2007).
- D. Marr, T. Poggio, "From understanding computation to understanding neural circuitry" (AIM-357, Massachusetts Institute of Technology 201 Vassar Street, W59-200, 1976), vol. 357.
- B. C. Love, Levels of biological plausibility. *Philos. Trans. R. Soc. B Biol. Sci.* **376**, 20190632 (2021).
- D. Marr, *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information* (MIT Press, 2010).
- J. Bickle, Marr and reductionism. *Top. Cogn. Sci.* **7**, 299–311 (2015).
- J. W. Krakauer, A. A. Ghazanfar, A. Gomez-Marín, M. A. McCliver, D. Poeppel, Neuroscience needs behavior: Correcting a reductionist bias. *Neuron* **93**, 480–490 (2017).
- D. Hebb, *The Organization of Behavior: A Neuropsychological Theory* (John Wiley and Sons Inc., 1949).
- G. Buzsáki, Neural syntax: Cell assemblies, synapses, and readers. *Neuron* **68**, 362–385 (2010).
- H. Eichenbaum, Barlow versus Hebb: When is it time to abandon the notion of feature detectors and adopt the cell assembly as the unit of cognition? *Neurosci. Lett.* **680**, 88–93 (2018).
- R. Quian Quiroga, L. Reddy, G. Kreiman, C. Koch, I. Fried, Invariant visual representation by single neurons in the human brain. *Nature* **435**, 1102–1107 (2005).
- K. D. Harris, J. Csicsvari, H. Hirase, G. Dragoi, G. Buzsáki, Organization of cell assemblies in the hippocampus. *Nature* **424**, 552–556 (2003).
- J. M. O'Keefe, L. Nadel, J. O'Keefe, *The Hippocampus as a Cognitive Map* (Clarendon Press, 1978).
- E. I. Moser, E. Kropff, M.-B. Moser, Place cells, grid cells, and the brain's spatial representation system. *Ann. Rev. Neurosci.* **31**, 69–89 (2008).
- A. Horner, J. Bisby, E. Zotow, D. Bush, N. Burgess, Grid-like processing of imagined navigation. *Curr. Biol.* **26**, 842–847 (2016).
- J. L. S. Bellmund, P. Gärdenfors, E. I. Moser, C. F. Doeller, Navigating cognition: Spatial codes for human thinking. *Science* **362**, eaat6766 (2018).
- T. E. J. Behrens, T. H. Muller, J. C. R. Whittington, S. Mark, A. B. Baram, K. L. Stachenfeld, Z. Kurth-Nelson, What is a cognitive map? Organizing knowledge for flexible behavior. *Neuron* **100**, 490–509 (2018).
- R. M. Mok, B. C. Love, A non-spatial account of place and grid cells based on clustering models of concept learning. *Nat. Commun.* **10**, 5685 (2019).
- E. C. Tolman, Cognitive maps in rats and men. *Psychol. Rev.* **55**, 189–208 (1948).
- B. C. Love, D. L. Medin, T. M. Gureckis, SUSTAIN: A network model of category learning. *Psychol. Rev.* **111**, 309–332 (2004).
- B. C. Love, T. M. Gureckis, Models in search of a brain. *Cogn. Affect. Behav. Neurosci.* **7**, 90–108 (2007).
- T. Davis, B. C. Love, A. R. Preston, Learning the exception to the rule: Model-based fMRI reveals specialized representations for surprising category members. *Cereb. Cortex* **22**, 260–273 (2012).
- M. L. Mack, B. C. Love, A. R. Preston, Dynamic updating of hippocampal object representations reflects new conceptual knowledge. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 13203–13208 (2016).
- M. L. Mack, A. R. Preston, B. C. Love, Ventromedial prefrontal cortex compression during concept learning. *Nat. Commun.* **11**, 46 (2020).
- M. B. Broschard, J. Kim, B. C. Love, E. A. Wasserman, J. H. Freeman, Prelimbic cortex maintains attention to category-relevant information and flexibly updates category representations. *Neurobiol. Learn. Mem.* **185**, 107524 (2021).
- C. W. Reynolds, in *Proceedings of the 14th Annual Conference on Computer Graphics and Interactive Techniques - SIGGRAPH '87* (ACM Press, 1987), pp. 25–34.
- R. Koster, M. J. Chadwick, Y. Chen, D. Berron, A. Banino, E. Düzel, D. Hassabis, D. Kumaran, Big-loop recurrence within the hippocampal system supports integration of information across episodes. *Neuron* **99**, 1342–1354.e6 (2018).
- A. Attardo, J. E. Fitzgerald, M. J. Schnitzer, Impermanence of dendritic spines in live adult CA1 hippocampus. *Nature* **523**, 592–596 (2015).
- M. L. Mack, B. C. Love, A. R. Preston, Building concepts one episode at a time: The hippocampus and concept formation. *Neurosci. Lett.* **680**, 31–38 (2018).
- R. N. Shepard, C. I. Hovland, H. M. Jenkins, Learning and memorization of classifications. *Psychol. Monogr. Gen. Appl.* **75**, 1–42 (1961).
- G. Šimić, I. Kostović, B. Winblad, N. Bogdanović, Volume and number of neurons of the human hippocampal formation in normal aging and Alzheimer's disease. *J. Comp. Neurol.* **379**, 482–494 (1997).
- N. Sukenik, O. Vinogradov, E. Weinreb, M. Segal, A. Levina, E. Moses, Neuronal circuits overcome imbalance in excitation and inhibition by adjusting connection numbers. *Proc. Natl. Acad. Sci. U.S.A.* **118**, e2018459118 (2021).
- H. Barlow, Redundancy reduction revisited. *Network* **12**, 241–253 (2001).
- N. S. Narayanan, Redundancy and synergy of neuronal ensembles in motor cortex. *J. Neurosci.* **25**, 4207–4216 (2005).
- J. L. Puchalla, E. Schneidman, R. A. Harris, M. J. Berry, Redundancy in the population code of the retina. *Neuron* **46**, 493–504 (2005).
- M. Jung, S. Wiener, B. McNaughton, Comparison of spatial firing characteristics of units in dorsal and ventral hippocampus of the rat. *J. Neurosci.* **14**, 7347–7356 (1994).
- J. Poppenk, H. R. Evensmoen, M. Moscovitch, L. Nadel, Long-axis specialization of the human hippocampus. *Trends Cogn. Sci.* **17**, 230–240 (2013).
- I. K. Brunec, B. Bellana, J. D. Ozubko, V. Man, J. Robin, Z. X. Liu, C. Grady, R. S. Rosenbaum, G. Winocur, M. D. Barense, M. Moscovitch, Multiple scales of representation along the hippocampal anteroposterior axis in humans. *Curr. Biol.* **28**, 2129–2135.e6 (2018).
- S. Tanni, W. de Corthi, C. Barry, State transitions in the statistically stable place cell population correspond to rate of perceptual change. *Curr. Biol.* **32**, 3505–3514.e7 (2022).
- J. S. Bowers, On the biological plausibility of grandmother cells: Implications for neural network theories in psychology and neuroscience. *Psychol. Rev.* **116**, 220–251 (2009).
- D. C. Plaut, J. L. McClelland, Locating object knowledge in the brain: Comment on Bowers's (2009) attempt to revive the grandmother cell hypothesis. *Psychol. Rev.* **117**, 284–288 (2010).
- R. Quian Quiroga, G. Kreiman, Measuring sparseness in the brain: Comment on Bowers (2009). *Psychol. Rev.* **117**, 291–297 (2010).
- J. Su'ević, A. C. Schapiro, A neural network model of hippocampal contributions to category learning. bioRxiv 2022.01.12.476051 [Preprint]. 13 January 2022. <https://doi.org/10.1101/2022.01.12.476051>.
- D. N. Barry, E. A. Maguire, Remote memory and the hippocampus: A constructive critique. *Trends Cogn. Sci.* **23**, 128–142 (2019).
- K. L. Spalding, O. Bergmann, K. Alkass, S. Bernard, M. Salehpour, H. B. Huttner, E. Boström, I. Westerlund, C. Vial, B. A. Buchholz, G. Possnert, D. C. Mash, H. Druid, J. Frisén, Dynamics of hippocampal neurogenesis in adult humans. *Cell* **153**, 1219–1227 (2013).
- L. Thompson, P. Best, Long-term stability of the place-field activity of single units recorded from the dorsal hippocampus of freely behaving rats. *Brain Res.* **509**, 299–308 (1990).
- H. S. Wirthshafter, J. F. Disterhoft, *In vivo* multi-day calcium imaging of CA1 hippocampus in freely moving rats reveals a high preponderance of place cells with consistent place fields. *J. Neurosci.* **42**, 4538–4554 (2022).
- A. Gilboa, M. Moscovitch, No consolidation without representation: Correspondence between neural and psychological representations in recent and remote memory. *Neuron* **109**, 2239–2255 (2021).

48. K. C. Bittner, A. D. Milstein, C. Grienberger, S. Romani, J. C. Magee, Behavioral time scale synaptic plasticity underlies CA1 place fields. *Science* **357**, 1033–1036 (2017).
49. J. C. Magee, C. Grienberger, Synaptic plasticity forms and functions. *Ann. Rev. Neurosci.* **43**, 95–117 (2020).
50. S. Zhao, X. Tang, W. Tian, S. Partarrieu, R. Liu, H. Shen, J. Lee, S. Guo, Z. Lin, J. Liu, Tracking neural activity from the same cells during the entire adult life of mice. *Nat. Neurosci.* **26**, 696–710 (2023).
51. C. MacDonald, K. Lepage, U. Eden, H. Eichenbaum, Hippocampal “time cells” bridge the gap in memory for discontinuous events. *Neuron* **71**, 737–749 (2011).
52. J. Gray, in *The Neuropsychology of Anxiety: An Enquiry into the Functions of the Septo-Hippocampal System* (Oxford University Press, 1982), vol. 5, pp. 469–484.
53. O. Vinogradova, Hippocampus as comparator: Role of the two input and two output systems of the hippocampus in selection and registration of information. *Hippocampus* **11**, 578–598 (2001).
54. J. E. Lisman, A. A. Grace, The hippocampal-VTA loop: Controlling the entry of information into long-term memory. *Neuron* **46**, 703–713 (2005).
55. N. Ketz, S. G. Morkonda, R. C. O'Reilly, Theta coordinated error-driven learning in the hippocampus. *PLOS Comput. Biol.* **9**, e1003067 (2013).
56. A. C. Schapiro, N. B. Turk-Browne, M. M. Botvinick, K. A. Norman, Complementary learning systems within the hippocampus: A neural network modelling approach to reconciling episodic memory with statistical learning. *Philos. Trans. R. Soc. B Biol. Sci.* **372**, 20160049 (2017).
57. Y. Zheng, X. L. Liu, S. Nishiyama, C. Ranganath, R. C. O'Reilly, Correcting the hebbian mistake: Toward a fully error-driven hippocampus. *PLOS Comput. Biol.* **18**, e1010589 (2022).
58. A. R. Preston, Y. Shrager, N. M. Dudukovic, J. D. Gabrieli, Hippocampal contribution to the novel use of relational information in declarative memory. *Hippocampus* **14**, 148–152 (2004).
59. A. Schapiro, L. Kustner, N. Turk-Browne, Shaping of object representations in the human medial temporal lobe based on temporal regularities. *Curr. Biol.* **22**, 1622–1627 (2012).
60. Y. Burak, I. R. Fiete, Accurate path integration in continuous attractor network models of grid cells. *PLOS Comput. Biol.* **5**, e1000291 (2009).
61. J. C. R. Whittington, T. H. Muller, S. Mark, G. Chen, C. Barry, N. Burgess, T. E. J. Behrens, The tolman-eichenbaum machine: Unifying space and relational memory through generalization in the hippocampal formation. *Cell* **183**, 1249–1263.e23 (2020).
62. D. George, R. V. Rikhye, N. Gothoskar, J. S. Guntupalli, A. Dedieu, M. Lázaro-Gredilla, Clone-structured graph representations enable flexible learning and vicarious evaluation of cognitive maps. *Nat. Commun.* **12**, 2392 (2021).
63. C. R. Bowman, T. Iwashita, D. Zeithamova, Tracking prototype and exemplar representations in the brain across learning. *eLife* **9**, e59360 (2020).
64. J. K. Kruschke, ALCOVE: An exemplar-based connectionist model of category learning. *Psychol. Rev.* **99**, 22–44 (1992).
65. R. M. Nosofsky, M. A. Gluck, T. J. Palmeri, S. C. Mckinley, P. Glauthier, Comparing modes of rule-based classification learning: A replication and extension of Shepard, Hovland, and Jenkins (1961). *Mem. Cognit.* **22**, 352–369 (1994).
66. A. Banino, C. Barry, B. Uria, C. Blundell, T. Lillicrap, P. Mirowski, A. Pritzel, M. J. Chadwick, T. Degris, J. Modayil, G. Wayne, H. Soyer, F. Viola, B. Zhang, R. Goroshin, N. Rabinowitz, R. Pascanu, C. Beattie, S. Petersen, A. Sadik, S. Gaffney, H. King, K. Kavukcuoglu, D. Hassabis, R. Hadsell, D. Kumaran, Vector-based navigation using grid-like representations in artificial agents. *Nature* **557**, 429–433 (2018).
67. F. Sargolini, M. Fyhn, T. Hafting, B. L. McNaughton, M. P. Witter, M. B. Moser, E. I. Moser, Conjunctive representation of position, direction, and velocity in entorhinal cortex. *Science* **312**, 758–762 (2006).

#### Acknowledgments

**Funding:** This work was supported by the Medical Research Council UK (MC\_UU\_00030/7) and a Leverhulme Trust Early Career Fellowship (Leverhulme Trust, Isaac Newton Trust: SUAI/053 G100773, SUAI/056 G105620, and ECF-2019-110) to R.M.M. and the Wellcome Trust (WT106931MA), ESRC (ES/W007347/1), and a Royal Society Wolfson Fellowship (18302) to B.C.L.

**Author contributions:** R.M.M.: Conceptualization, data curation, formal analysis, funding acquisition, methodology, software, visualization, writing—original draft, and writing—review and editing. B.C.L.: Conceptualization, formal analysis, funding acquisition, methodology, resources, supervision, visualization, writing—original draft, and writing—review and editing.

**Competing interests:** The authors declare that they have no competing interests. **Data and materials availability:** Code and simulation results are available at GitHub (<https://github.com/robmok/multiunit-cluster>), Zenodo (<https://zenodo.org/badge/latestdoi/308407059>), and OSF (<https://osf.io/uf8pa/>). All data needed to evaluate the conclusions in the paper are present in the paper and/or the Supplementary Materials. For the purpose of open access, the author has applied a Creative Commons Attribution (CC BY) licence to any author-accepted manuscript version arising from this submission. Fruit images were obtained from vecteezy.com.

Submitted 1 September 2022

Accepted 20 June 2023

Published 21 July 2023

10.1126/sciadv.ade6903

## **A multilevel account of hippocampal function in spatial and concept learning: Bridging models of behavior and neural assemblies**

Robert M. Mok and Bradley C. Love

*Sci. Adv.*, **9** (29), eade6903.  
DOI: 10.1126/sciadv.ade6903

### **View the article online**

<https://www.science.org/doi/10.1126/sciadv.ade6903>

### **Permissions**

<https://www.science.org/help/reprints-and-permissions>

Use of this article is subject to the [Terms of service](#)

---

*Science Advances* (ISSN ) is published by the American Association for the Advancement of Science. 1200 New York Avenue NW, Washington, DC 20005. The title *Science Advances* is a registered trademark of AAAS.  
Copyright © 2023 The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original U.S. Government Works. Distributed under a Creative Commons Attribution License 4.0 (CC BY).