

# Generalisation of structural knowledge in hippocampal – prefrontal circuits

**Veronika Samborska\*** ([veronika.samborska@ndcn.ox.ac.uk](mailto:veronika.samborska@ndcn.ox.ac.uk))

Nuffield Department of Clinical Neurosciences, University of Oxford, 13 Mansfield Rd, Oxford OX1 3SR, UK

**Thomas Akam\*** ([thomas.akam@psy.ox.ac.uk](mailto:thomas.akam@psy.ox.ac.uk))

Department of Experimental Psychology, University of Oxford, 13 Mansfield Rd, Oxford OX1 3SR, UK

**James L. Butler** ([ucbtjbu@ucl.ac.uk](mailto:ucbtjbu@ucl.ac.uk))

Institute of Neurology, UCL, Queen Square, London WC1N 3BG, UK

**Mark E. Walton†** ([mark.walton@psy.ox.ac.uk](mailto:mark.walton@psy.ox.ac.uk))

Department of Experimental Psychology, University of Oxford, 13 Mansfield Rd, Oxford OX1 3SR, UK

**Timothy E. Behrens†** ([behrens@fmrib.ox.ac.uk](mailto:behrens@fmrib.ox.ac.uk))

Nuffield Department of Clinical Neurosciences, University of Oxford, John Radcliffe Hospital, Oxford OX3 9DU, UK

\*, † Equal contribution.

## Abstract:

The ability to generalise previously learned knowledge to solve novel analogous problems relies on formation of representations that are abstracted from sensory states. Little is known about how the brain generalises abstract representations while maintaining the content of individual experiences. Here we present a novel behavioural paradigm for investigating generalisation of structural knowledge in mice and report electrophysiological findings from single neurons in hippocampus and prefrontal cortex. Mice serially performed a set of reversal learning tasks, which shared the same structure (e.g., one choice port is good at a time), but had different physical configurations and hence different sensory and motor representations. Subjects' performance on novel configurations improved with the number of configurations they had already learned, demonstrating generalisation of knowledge. As in spatial remapping experiments, many hippocampal neurons responded differently in different configurations – here tasks rather than spatial environments. In contrast, prefrontal representations were more general and reflected different stages of the trial irrespective of the current physical configuration. Population analyses showed that although the structure of each task was represented strongly in both regions, different hippocampal neuronal assemblies participated in each task's representation. In contrast, neuronal patterns in PFC generalised between different configurations.

**Keywords:** structure learning; hippocampus; prefrontal cortex; value-based decision making

## Background

Recent progress in cognitive neuroscience has provided us with a formal understanding of how our brain learns from direct experience (Cohen et al. 2012). A major open challenge is to understand the broader

class of behaviours where prior knowledge is generalised to solve new problems.

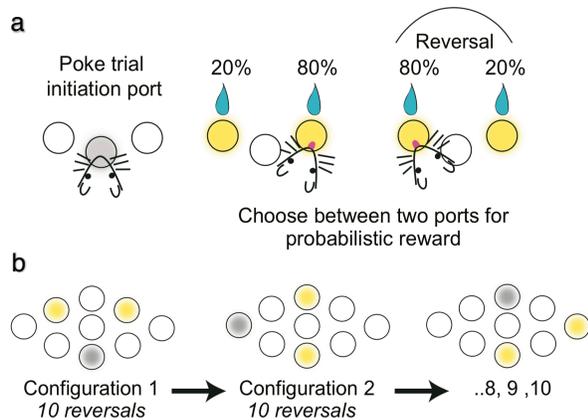
The ability to make appropriate inferences that go far beyond one's experience is thought to rely on the brain's internal model of the world, termed a cognitive map (Tolman, 1948). Cognitive maps are well studied in spatial domains, where we have detailed knowledge of underlying cellular codes in the hippocampal formation (Grieves & Jeffery, 2017). Recent data from rodents (Aronov, Nevers & Tank, 2017) and humans (Constantinescu, O'Reilly & Behrens, 2016) suggests that the same cellular mechanisms might encode complex relationships that organise knowledge outside the spatial domain. Importantly, map-like representations of non-spatial models resembling grid codes have been recorded in human fMRI in prefrontal cortex (Constantinescu et al. 2016). Here we combined a novel behavioural paradigm in mice with silicon probe recordings, to investigate how hippocampus and prefrontal cortex represent structure knowledge to enable generalisation of structure learning in a non-spatial domain.

## Methods

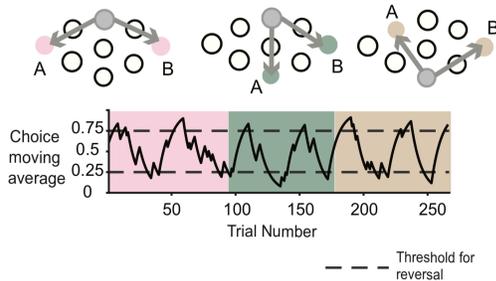
Mice were trained to solve a probabilistic reversal-learning paradigm where they initiated a trial by poking in one port, then chose between two other ports for a probabilistic reward (Fig 1a). Once the subjects consistently chose the high reward probability port, the reward contingencies reversed. When subjects had competed ten reversals on a given port configuration (termed a 'task'), they were moved onto another task with a different port configuration (Fig 1b). All of the



tasks shared the same trial structure (initiate → choose) and a common abstract rule (one port has high and one low reward probability, with occasional reversals), but the specific location of the ports and hence the actions required to perform trials were different in each configuration. We used silicon probes to record from hippocampal CA1 (396 neurons,  $n = 4$  mice) and medial prefrontal cortex (567 neurons,  $n = 4$  mice). In recording sessions mice typically performed at least 4 reversals in each of 3 different configurations (Fig 2).



**Figure 1: Trial structure of the probabilistic reversal-learning paradigm.** **a)** Mice poked in an initiation port then chose between two choice ports for a probabilistic reward. Reward contingencies reversed after the animal consistently chose the ‘good’ port. **b)** Example configurations used in each task. Subjects completed 10 reversals on each layout before moving to a new poke port layout.



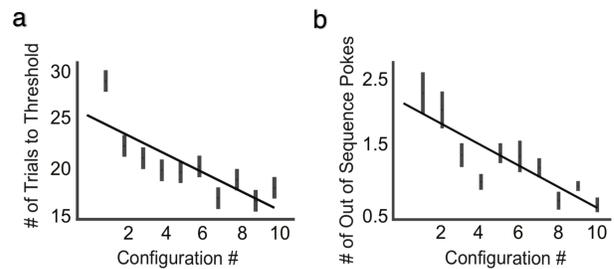
**Figure 2: Example of the configurations used in a single recording session.** Mice typically completed three configurations each consisting of four reversals in each recording session.

## Results

### Behavioural Results

Mice got better at tracking the good port over the course of each physical configuration of the task (i.e. fewer trials to criterion), but critically also showed

improvement across tasks with different configurations (Fig 3a). Moreover, they also got better at following within trial structure (initiate → choose) across tasks, making fewer pokes to invalid ports (Fig 3b). This demonstrates generalisation of learning and suggests that mice may have developed sensory invariant/abstract representation of the structure of the task.

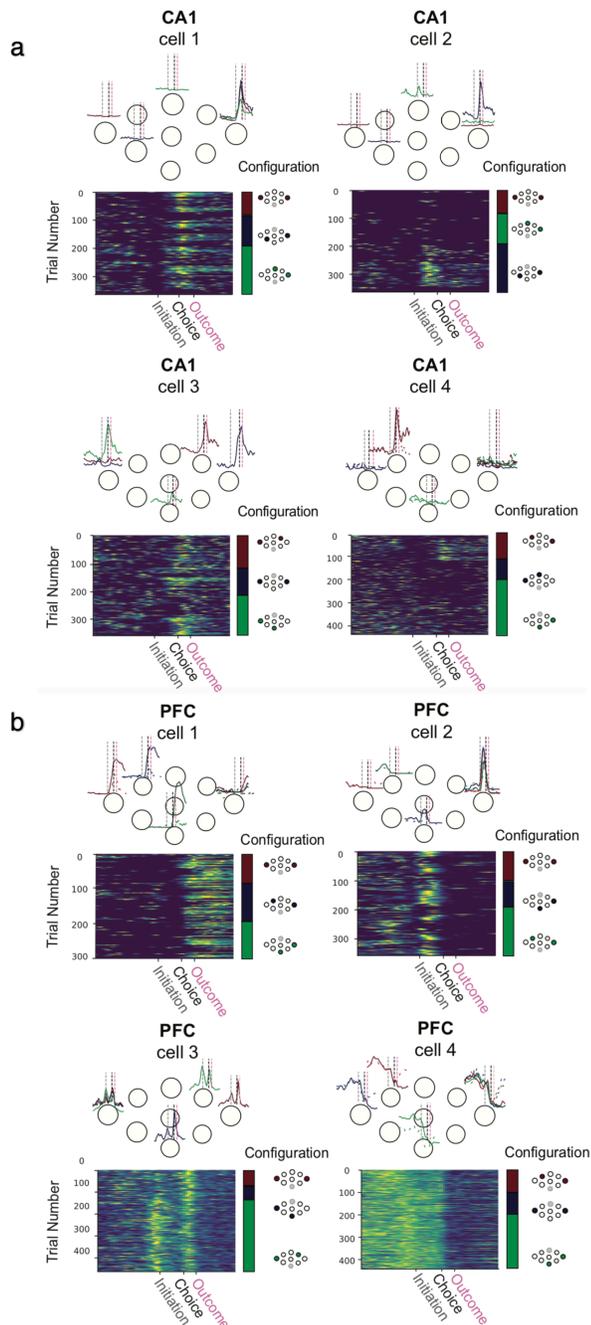


**Figure 3: Generalisation of structural knowledge in mice.**

**a)** Median number of trials mice took to reach the reversal threshold on each task configuration of a reversal-learning problem. **b)** Median number of pokes per trial mice made to a choice port that was no longer available because they already chose the other choice port.

### Single Unit Results

A substantial proportion of hippocampal neurons fired selectively when mice entered a particular physical port (Fig 4a, cell 1). However, when the task configuration changed, the activity of many hippocampal units ‘remapped’. Some units fired at a given port in one configuration but not when the same port was visited in a different configuration (Fig 4a, cell 2). Other units fired selectively to one of the choice ports in each configuration (Fig 4a, cell 3). Many cells had conjunctive place x reward coding, modulating their firing to a particular poke port based on whether it was rewarded or not (Fig 4a, cell 4). In contrast, prefrontal representations appeared to be more invariant across configurations. Many cells fired selectively when one of the choice ports was rewarded irrespective of its location in all three tasks (Fig 4b, cell 1). We also found units that fired at the shared port in all tasks irrespective of whether it was rewarded or not (Fig 4b, cell 2) and neurons that had multiple peaks throughout the trial (e.g., selective for both initiation and choice states) irrespective of the task layout (Fig 4b, cell 1). Furthermore, many neurons in PFC generalised their complex temporal trial dynamics across tasks (Fig 4b, cell 4). PFC cells therefore appeared to represent task states irrespective of the physical locations of the ports.



**Figure 4: Example neurons from CA1 and PFC.** Upper panels show normalised mean firing rates aligned to choice port entry time. Colours indicate configuration (green, blue, red). Vertical dashed lines indicate initiation (grey), choice (black) and outcome (pink) times. CA1 cell 4, PFC cells 1,2 and 4 were split by rewarded and non-rewarded trials. Solid lines indicate mean firing rates on rewarded trials, dashed lines are rates on non-rewarded trials. Lower panels show heat maps of normalised firing rate as a function of time within trial and trial number. Configurations are indicated on the right. **a) ‘Place cell’ like firing and ‘remapping’ in**

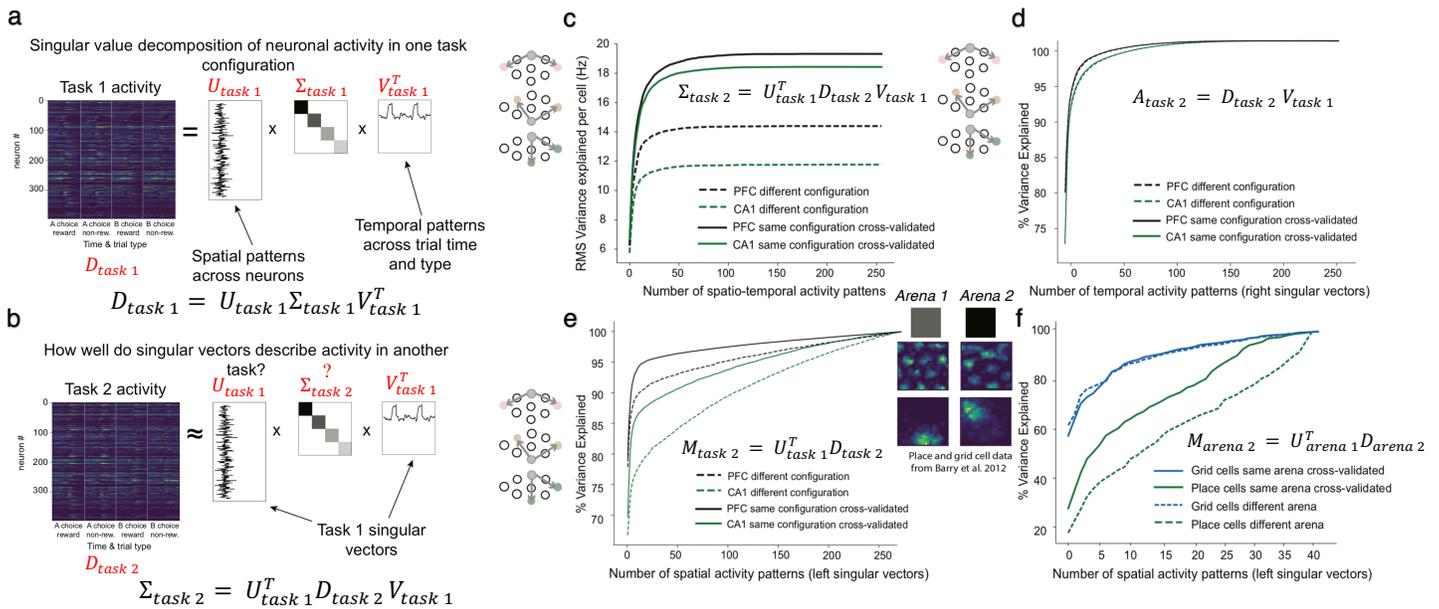
**hippocampus.** Cell 1 is a spatial cell that fired at a particular poke across all three configurations. Cell 2 is a ‘remapping’ cell that only fired at the far-right port in the context of one of the configurations. Cell 3 fired at different choice ports in different configurations. Cell 4 fired for one choice port and continued to fire following the outcome if the choice was rewarded. **b) Sensory invariant representations of trial structure in prefrontal cortex.** Cell 1 was active at one of the choice ports when it was rewarded in all configurations. Cell 2 fired at the shared choice port in all configurations. Cell 3 was active at both initiation and choice ports across all configurations. Cell 4 fired prior to all choice ports in all configurations, decreased its firing at choice, and only increased its firing rate again if the choice was not rewarded.

## Population Results

Using singular value decomposition (SVD), we asked whether population activity in PFC and CA1 shared the same low dimensional space across tasks (Fig 5a).

We found that low dimensional spatio-temporal patterns generalised better (i.e. explained more variance in a task with a different configuration) in PFC than CA1 (Fig 5c). Next, to look at how well the structure of each task is represented in each region we removed the constraint for the temporal activity patterns to be paired with particular spatial activity patterns and looked at how well we could explain activity in a task using just the temporal patterns from a task with a different port configuration (Fig 5d). Singular vectors corresponding to temporal patterns across time and trial type generalised near perfectly across tasks, confirming that in both regions temporal patterns that described activity in one task described activity in all tasks.

Finally, to look at how well the correlational structure between neurons was preserved across different configurations we removed the constraint for the spatial activity patterns to be paired with particular temporal activity patterns. Singular vectors corresponding to just the spatial patterns of activity across neurons generalised better in PFC than CA1, confirming that CA1 ‘remapped’ more than PFC between tasks (Fig 5e). This suggests that even though both regions had task representations that generalised across configurations, different neuronal assemblies participated in each task’s representation in CA1 but less so in PFC. We performed the same analysis on hippocampal place cells and entorhinal grid cells recorded in different physical environments by Barry et al. (2012). Analogous to our PFC cells, entorhinal grid cells maintained their relative firing positions (i.e. their correlation structure) across environments, while CA1 place cells ‘remapped’ (Fig 5f).



**Figure 5: Population Analyses Using Singular Value Decomposition.** **a)** SVD was first used to decompose a data matrix comprising the activity of each neuron across time and trial types from one task into the product of three matrices which linked a set of temporal patterns across trial type and time (rows of  $V_{task 1}^T$ ) to a set of spatial patterns across neurons (columns of  $U_{task 1}$ ). **b)** The sets of temporal and spatial patterns from the first task  $U_{task 1}$  and  $V_{task 1}^T$  were used to find the strength of the links  $\Sigma_{task 2}$  between these spatial and temporal vectors in a new task with a different physical configuration of the ports  $D_{task 2}$ . In all plots solid lines represent cumulative root mean square or % of the variance explained using singular vectors from the **same** task as the data matrix. Dashed lines represent cumulative root mean square or % of the variance explained using singular vectors from one task to explain a data matrix from a **different** task configuration. **c) Spatio-temporal activity patterns generalised better across tasks in PFC than CA1.** Lines represent cumulative singular values along the diagonal of  $\Sigma$ . **d) Temporal activity patterns generalise perfectly across tasks in both regions.** Lines represent cumulative sums of squares of values along the columns of  $A_{task 1/2}$  as a result of projecting the singular vectors corresponding to temporal patterns across time and trial onto the data matrix. **e) Spatial activity patterns across neurons generalise well in PFC but less well in hippocampus analogous to grid cells in entorhinal cortex and place cells in CA1 in two different spatial contexts (f).** Lines represent cumulative sums of squares of values along the columns of  $M_{task 1/2}$  in **(e)** or  $M_{arena 1/2}$  in **(f)** as a result of projecting the data matrix from one of the tasks **(e)** or heatmap from one of the contexts **(f)** onto the singular vectors of the correlational patterns across neurons in PFC and CA1 **(e)** and entorhinal grid cells and CA1 place cells **(f)**.

Our results provide preliminary evidence for common neuronal mechanisms underlying generalisation of structure knowledge in spatial and non-spatial domains.

## References

- Aronov, D., Nevers, R., & Tank, D. W. (2017). Mapping of a non-spatial dimension by the hippocampal-entorhinal circuit. *Nature*, 543 (7647), 719–722.
- Barry, C., Ginzberg, L. L., O’Keefe, J., & Burgess, N. (2012). Grid cell firing patterns signal environmental novelty by expansion. *Proceedings of the National Academy of Sciences*.
- Constantinescu, A. O., O’Reilly, J. X., & Behrens, T. E. J. (2016). Organizing conceptual knowledge in humans with a gridlike code. *Science*, 352 (6292), 1464–1468.
- Cohen, J. Y., Haesler, S., Vong, L., Lowell, B. B., & Uchida, N. (2012). Neuron-type-specific signals for reward and punishment in the ventral tegmental area. *Nature*.
- Grieves, R. M., & Jeffery, K. J. (2017). The representation of space in the brain. *Behavioural Processes*.
- Tolman, E. C. (1948). Cognitive maps in rats and men. *Psychological Review*.