

# The Helmholtz Hippocampus: A biologically plausible generative model of the hippocampal formation

Tom George<sup>1</sup>, Caswell Barry<sup>1,2</sup>, Kimberly Stachenfeld<sup>3,4</sup>, Claudia Clopath<sup>1,5</sup>, Tomoki Fukai<sup>6</sup>

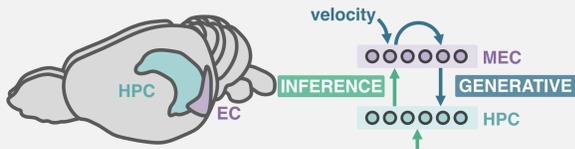


<sup>1</sup>Sainsbury Wellcome Centre, UCL, London <sup>2</sup>Department of Cell & Developmental Biology, UCL, London <sup>3</sup>Columbia University, New York <sup>4</sup>Google DeepMind, London UK <sup>5</sup>Imperial College, London <sup>6</sup>Okinawa Institute of Science and Technology, Japan

## 1 THE DUAL ROLE OF THE HIPPOCAMPUS

The hippocampal formation (HPC & MEC) has two main roles in navigation:

- 1. INFER self-location**  
place cells, grid cells etc.
- 2. GENERATE trajectories offline**  
replay<sup>[1]</sup>, memory consolidation<sup>[2]</sup>, planning<sup>[3]</sup>



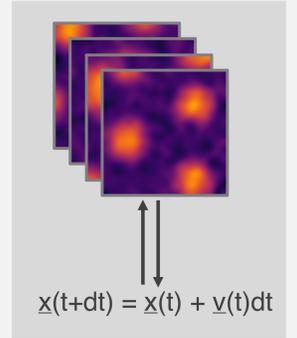
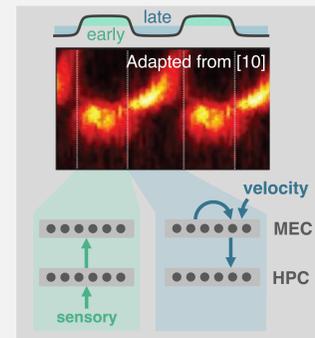
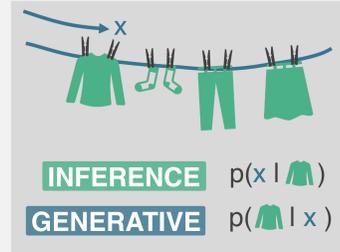
Existing models rely on biologically implausible assumptions<sup>[4,5,6]</sup>.

- Can a **biologically plausible** model account for these **inferential** and **generative** capacities.
- What would the **architecture, dynamics** and **learning rules** look like?
- What might it teach us about generative models in the brain?

TL;DR?  
See box 6

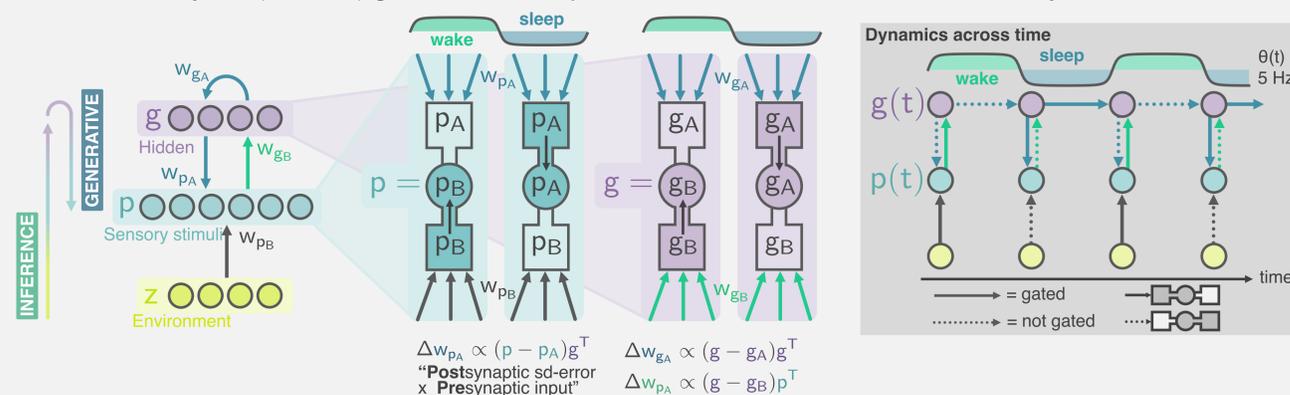
## 2 BACKGROUND

- MEC learns a sensory-agnostic map of space. HPC binds sensory stimuli onto positions in this space<sup>[7,5]</sup>.
- During behaviour, **inference** and **generation** occur in distinct **early** and **late** phase of the 5-10 Hz theta rhythm<sup>[8,9]</sup>.
- Path integration – a **“generative”** process – is a fundamental pillar of MEC function<sup>[11,4]</sup>.



## 3 THE HELMHOLTZ HIPPOCAMPUS

- A hierarchical network receives a continuous stream of sensory inputs.
- Inference** and **generative** pathways arrive at distinct **basal** and **apical** dendritic compartments.
- LFP theta-rhythm (5-10 Hz) gates which compartment drives the soma, thus which way information flows.



- Inference** and **generative sampling** occur in alternating 5-10 Hz **wake-sleep** phases.
- Hebbian learning**: minimizes the somatic “prediction errors”<sup>[12]</sup> until  $p_B \approx p_A$  and  $g_B \approx g_A$

## 4 THEORETICAL INTERPRETATION

Helmholtz machines<sup>[13]</sup> (aka Boltzmann machines) learn latent models by matching **inference** and **generative** distributions in alternating learning phases.

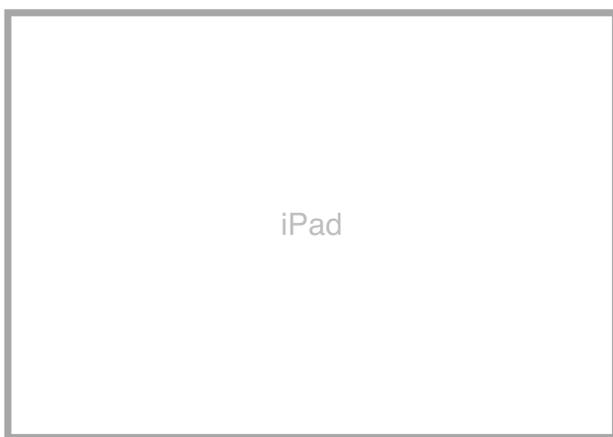
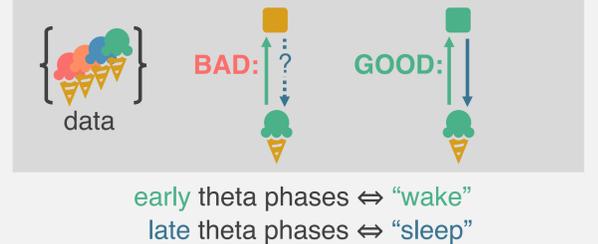
$$p_{\text{inf}}(p, g) = p_w(g|p)p(p)$$

$$p_{\text{gen}}(p, g) = p_w(p|g)p_w(g)$$

$$\mathcal{L}(w) = D_{\text{KL}}(p_{\text{inf}} || p_{\text{gen}})$$

...our **Hebbian learning rule** is derived from this **loss function**<sup>[14]</sup>.

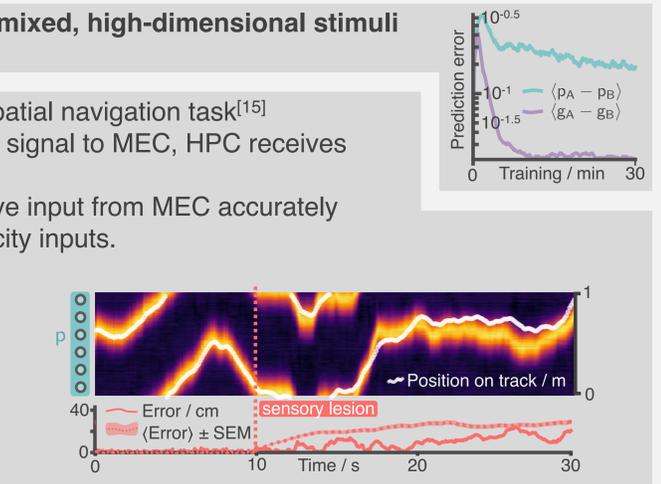
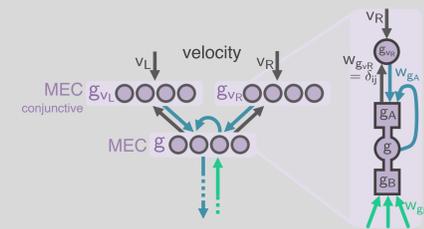
“What I cannot **[generate]** I do not **understand**”



## 5 RESULTS

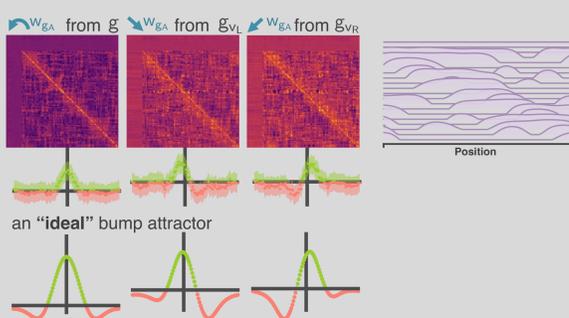
**1. COMPRESSION:** HPC learns to efficiently **compress mixed, high-dimensional stimuli** deriving from a small number of independent latents.

- Initially un-tuned conjunctive inputs provide velocity signal to MEC, HPC receives spatially-selective “place cell” inputs.
- When sensory input is lesioned, top-down generative input from MEC accurately maintains position estimate by “integrating” its velocity inputs.



**3. RING-ATTRACTOR:** Weights in the hidden layer reveal a stable calibrated ring-attractor.

This theoretically justified<sup>[16]</sup> solution is hard to achieve with local learning rules<sup>[17,18]</sup>.



**4. REMAPPING AND TRANSFER LEARNING:** We tested “remapping” by **fixing the MEC weights** and **shuffling sensory input**

- Model quickly **remastered new environment** by transferring (not relearning) path integration.
- MEC hidden representations reformed with constant phase shift, reminiscent of real grid cells<sup>[14]</sup>.



## 6 CONCLUSIONS

- Hippocampal **architecture, dynamics** and **learning rules** are highly similar to those of a **Helmholtz machine**. HPC maps sensory input to MEC which learns hidden latent structure.
- LFP Theta oscillations** rapidly gate information flow to generate **“wake-sleep”** learning phases with exclusively **local learning rules**.
- Various navigational function, including self-location, path-integration and transfer learning, can thence be explained **without backpropagation**.

NeurIPS 2023 paper:



[1] Lee & Wilson (2002), Memory of sequential experience in the hippocampus during slow wave sleep [2] Carr et al. (2011), Hippocampal replay in the awake state: a potential substrate for memory consolidation and retrieval [3] Spiers and Maguire (2006), Thoughts, behaviour, and brain dynamics during navigation in the real world [4] Banino et al. (2018), Vector-based navigation using grid-like representations in artificial agents [5] Whittington et al. (2020), The Tolman-Eichenbaum Machine: Unifying Space and Relational Memory through Generalization in the Hippocampal Formation [6] Sorscher et al. (2022), A unified theory for the computational and mechanistic origins of grid cells [7] Bellmund et al. (2018), Navigating cognition: Spatial codes for human thinking [8] Hasselmo et al. (2002), A Proposed Function for Hippocampal Theta Rhythm: Separate Phases of Encoding and Retrieval Enhance Reversal of Prior Learning [9] Sanders and Lisman (2015), Grid Cells and Place Cells: An Integrated View of their Navigational and Memory Function [10] Wang et al. (2020), Alternating sequences of future and past behavior encoded within hippocampal theta oscillations [11] Dorrell et al. (2023), Actionable Neural Representations: Grid Cells from Minimal Constraints [12] Urbanczik and Senn (2014), Learning by the Dendritic Prediction of Somatic Spiking [13] Dayan et al. (1995), The Helmholtz Machine [14] Bredenberg et al. (2021), Impression learning: Online representation learning with synaptic plasticity [15] George et al. (2024), RatinABox, a toolkit for modelling locomotion and neuronal activity in continuous environments [16] Zhang (1996), Representation of spatial orientation by the intrinsic dynamics of the head-direction cell ensemble: a theory [17] Burak and Fiete (2009), Accurate path integration in continuous attractor network models of grid cells [18] Vafidis et al. (2022), Learning accurate path integration in ring attractor models of the head direction system [19] Fyhn et al. (2007), Hippocampal remapping and grid realignment in entorhinal cortex.