

# Goal-seeking compresses neural codes for space in the human hippocampus and orbitofrontal cortex

## Highlights

- Humans achieved context-dependent navigation, taking multiple steps to reach a final goal
- Goals induced compression of spatial maps in the hippocampus and orbitofrontal cortices
- Compression tracked performance and emerged only with full knowledge of goal locations
- Captured by a place cell model that jointly encodes current and prospective locations

## Authors

Paul S. Muhle-Karbe, Hannah Sheahan, Giovanni Pezzulo, Hugo J. Spiers, Samson Chien, Nicolas W. Schuck, Christopher Summerfield

## Correspondence

p.muhle-karbe@bham.ac.uk (P.S.M.-K.),  
christopher.summerfield@psy.ox.ac.uk (C.S.)

## In brief

How do goals affect the representation of space? Muhle-Karbe et al. show that spatial maps in the hippocampus and orbitofrontal cortex are “compressed” during goal-directed navigation with goals that share a prospective pathway clustering together in representational space, as if participants were simultaneously imagining themselves at current and goal locations.

Article

# Goal-seeking compresses neural codes for space in the human hippocampus and orbitofrontal cortex

Paul S. Muhle-Karbe,<sup>1,2,3,10,12,\*</sup> Hannah Sheahan,<sup>1,4,10</sup> Giovanni Pezzulo,<sup>5</sup> Hugo J. Spiers,<sup>6</sup> Samson Chien,<sup>7</sup> Nicolas W. Schuck,<sup>7,8,9,11</sup> and Christopher Summerfield<sup>1,3,11,\*</sup>

<sup>1</sup>Department of Experimental Psychology, University of Oxford, Oxford OX2 6GG, UK

<sup>2</sup>School of Psychology, University of Birmingham, Birmingham B15 2SA, UK

<sup>3</sup>Centre for Human Brain Health, University of Birmingham, Birmingham B15 2SA, UK

<sup>4</sup>Google DeepMind, London EC4A 3TW, UK

<sup>5</sup>Institute of Cognitive Sciences and Technologies, National Research Council, 00185 Rome, Italy

<sup>6</sup>Department of Experimental Psychology, University College London, London WC1E 6BT, UK

<sup>7</sup>Max Planck Research Group NeuroCode, Max Planck Institute for Human Development, 14195 Berlin, Germany

<sup>8</sup>Max Planck UCL Centre for Computational Psychiatry and Aging Research, 14195 Berlin, Germany

<sup>9</sup>Institute of Psychology, Universität Hamburg, 20146 Hamburg, Germany

<sup>10</sup>These authors contributed equally

<sup>11</sup>Senior author

<sup>12</sup>Lead contact

\*Correspondence: [p.muhle-karbe@bham.ac.uk](mailto:p.muhle-karbe@bham.ac.uk) (P.S.M.-K.), [christopher.summerfield@psy.ox.ac.uk](mailto:christopher.summerfield@psy.ox.ac.uk) (C.S.)

<https://doi.org/10.1016/j.neuron.2023.08.021>

## SUMMARY

Humans can navigate flexibly to meet their goals. Here, we asked how the neural representation of allocentric space is distorted by goal-directed behavior. Participants navigated an agent to two successive goal locations in a grid world environment comprising four interlinked rooms, with a contextual cue indicating the conditional dependence of one goal location on another. Examining the neural geometry by which room and context were encoded in fMRI signals, we found that map-like representations of the environment emerged in both hippocampus and neocortex. Cognitive maps in hippocampus and orbitofrontal cortices were compressed so that locations cued as goals were coded together in neural state space, and these distortions predicted successful learning. This effect was captured by a computational model in which current and prospective locations are jointly encoded in a place code, providing a theory of how goals warp the neural representation of space in macroscopic neural signals.

## INTRODUCTION

Humans and other primates can use context to guide their decisions. During *instantaneous* choices, where stimuli evoke independent action-outcome mappings, contextual cues modulate neural encoding of information in sensory neocortex<sup>1–3</sup> and higher regions such as prefrontal cortex (PFC).<sup>4–8</sup> However, context can also influence *sequential* choices, where outcomes depend on a series of transitions between states and actions, e.g., when navigating to a spatial goal. Many species can navigate flexibly to distinct goals based on contextual information that is unobservable or maintained in memory.<sup>9,10</sup> For example, a person might find their way to the local hairdresser or post office depending on the purpose of an errand. Flexible, context-dependent navigation requires that space and goals are encoded in ways that avoid mutual interference, allowing the correct destination to be reached given the context. Here, we studied the neural and computational mechanisms that make this possible in humans.

Recordings from the rodent hippocampus (HC) and connected structures have revealed much about the neural representation of allocentric space.<sup>11</sup> In the HC, “place cells” code for the animal’s current location via spatially localized firing patterns called “place fields.” These collectively form an internal map of the local environment, with each cell firing at a slightly different spatial location.<sup>12</sup> Neural codes for space have also been identified using single-cell recordings in other species, including humans,<sup>13,14</sup> and gross spatial location can be read out from fMRI signals.<sup>15–18</sup> In rodents, changes in context can lead place cells to form new fields in different locations, a phenomenon known as remapping. Global (and partial) remapping has been observed after physical changes to the local environment, such as the introduction of novel colors, textures, or odors<sup>19</sup> or the repositioning of the testing apparatus in a new room.<sup>20</sup> However, remapping can also occur when the context is denoted by an unobservable variable, such as a latent task rule,<sup>21</sup> a noisy inference about the environment-generating process,<sup>22</sup> or a prospective pathway or destination.<sup>23–25</sup>

Sometimes, remapping may occur along a single dimension aligned with the gain of neural activity, which is called rate remapping.<sup>26</sup>

The context provided by a spatial goal can also distort the representation of space without provoking gross changes in the neural code. Place cells tend to over-represent behaviorally significant spatial locations and can accumulate around<sup>27–29</sup> or fire excess spikes at a rewarded position.<sup>30</sup> Place cells may also encode information about prospective as well as current locations on the spatial trajectory,<sup>31–33</sup> and information about future states has also been observed in hippocampal blood-oxygen-level-dependent (BOLD) signals<sup>34,35</sup> and intracranial recordings.<sup>13,34,36</sup> It has also been claimed that spatial goals may be directly coded in the hippocampal formation.<sup>9,37</sup> One recent study has reported a small but dedicated population of CA1 neurons whose activity covaries with the location of a rewarding stimulus. When changes to the environment cause global remapping, these cells show the same preserved activity pattern linked to reward proximity.<sup>38</sup> These could be “goal cells,” a putative class of neuron that codes directly for a location that the animal seeks to reach, rather than its current spatial location.<sup>39</sup> Recent recordings from the rodent orbitofrontal cortex (OFC) provide parallel evidence of goal coding where the future goal location is present in the neural population before the animal begins its navigation to the goal.<sup>40</sup>

Decoding a mixture of place and goal locations could produce a spatial representation that is warped by the animal’s intended destination, with regions of space containing two prospective goal locations being coded with a more similar neural code, and thus appearing closer together in the internal spatial map (we call this “goal-based spatial compression”). Here, we report that although different neocortical areas encode goals and locations in heterogeneous ways, strong goal-based spatial compression is observed in the BOLD signal recorded from the human HC and orbitofrontal cortices.

## RESULTS

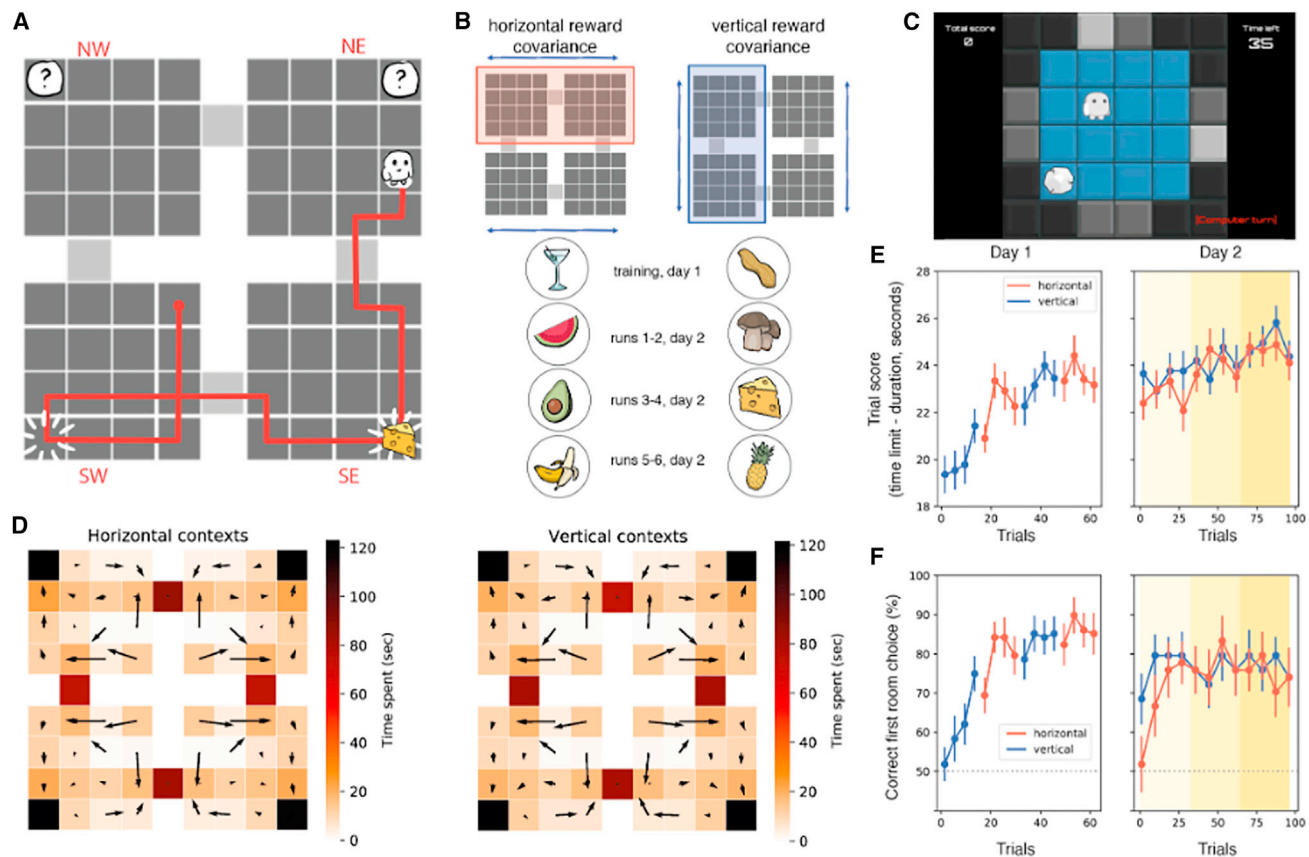
Human participants ( $n = 27$ ) performed a spatial navigation task that involved controlling an avatar that moved through a partially observable grid world composed of four discrete and interconnected rooms (Figure 1A). Participants saw a birds’ eye view of the currently occupied room (a single  $4 \times 4$  grid of squares); other rooms were not visible (Figure 1C). One grid square in each room contained a boulder, and across the entire environment, rewards were hidden under two of the four boulders. On each trial, the avatar spawned in a random room, and participants’ task was to move it (using buttons for up, down, left, and right) to collide with the two boulders that yielded rewards (within the minimum possible number of steps and in any order), avoiding those that were empty. Successful trial completion required both goals to be visited within a fixed time period.

Trials began with a contextual cue, which was a picture of one of the two food items (Figure 1B). Unbeknownst to participants, each cue revealed one reward location conditional on the other: half of the cues (“cue H”) indicated that the rewards were in rooms lying in the same horizontal axis, whereas the other half (“cue V”) indicated that the rewards were in rooms lying in the

same vertical axis (with neither disclosing which specific rooms where rewards could be found). Interleaving cues from trial to trial ensured that participants with perfect knowledge of the rules would on average display the same room occupancy probabilities across contexts. However, to ensure that participants also visited each room within each context, we introduced a “robot control” phase in each trial in which participants relinquished control of the avatar to a game controller, typically moving it to a suboptimal location. This manipulation allowed us to measure BOLD signals from locations that were off the shortest path taken by expert players<sup>41</sup> and ensured that room occupancy probabilities and transitions were well balanced across the experiment (Figure 1D).

After learning the task on an earlier day (see STAR Methods; Figures 1E and 1F, left panels), participants performed 96 trials across 6 scanner runs. In the scanner, we used three sets of physically distinct cues (food items) on runs 1–2, 3–4, and 5–6, respectively, requiring participants to generalize their knowledge of task structure across these three phases (Figure 1B). We plot behavioral results in Figure 1F. Although participants had no way of knowing whether or not the first boulder they encountered was rewarded, a participant with knowledge of the task structure can use this outcome in combination with the cue identity to exit the first room in the correct direction. Accordingly, on day 2, participants explored the start boulder on 98% of trials, and their first-choice accuracy increased rapidly across the first two runs, stabilizing at about 75% (lower panels), and was significantly above chance overall ( $t_{26} = 8.83$ ,  $p < 0.001$ ). Time taken to complete each (correct) trial continued to decrease across the experiment (Figure 1E, upper panels). At the end of the scanning session, participants completed a short quiz in which they were asked which room/s contained reward/s, given the presence or absence of rewards in other rooms. For example, “*You have just found a cheese in the top right room. Which room will the other cheese be in?*”. The mean score across participants was  $69\% \pm 12\%$  (chance performance was 34%). Quiz score positively correlated with both the average trial score on day 2 ( $r = 0.445$ ,  $p = 0.020$ ), and average first-choice accuracy on day 2 ( $r = 0.813$ ,  $p < 0.001$ ), suggesting that sequential decisions were guided by explicit knowledge about the task’s latent reward covariance structure.

To formulate neural predictions, we built a computational model that encoded the avatar’s location via simulated Gaussian place fields tiling an internal representation of the four rooms environment (Figure 2A). We read out the responses elicited across the neural population as each participant moved the avatar through the four rooms, by providing empirically observed trajectories from yoked human behavioral data as inputs to the model (Figure 2B). This allowed us to compute simulated representational dissimilarity matrices (RDMs) for each room and context (8 conditions), which were averaged before multidimensional scaling (MDS) was used to visualize their neural geometry. Without further elaboration, this model simply encodes the locations of the four rooms (Northeast [NE]; Northwest [NW], Southeast [SE], and Southwest [SW]) at the vertices of two perfectly parallel and aligned spatial maps, each visualized as a square plane denoting one context (Figure 2C, upper left panel). This unbiased geometry was obtained with the observed human



**Figure 1. Task structure and performance**

(A) Illustration of the four rooms environment and example reward locations under the vertical context. An example participant trajectory is shown overlaid in red. In this example, the agent starts in the Southwest (SW) room, explores the start boulder and does not find a cheese reward. As cheese rewards covary vertically, the two cheese rewards must therefore be in the SE and NE rooms.

(B) Different contexts signaled different reward covariances: on day 1 (training) martini rewards appeared in vertically adjacent rooms, while peanut rewards appeared in horizontally adjacent rooms. On day 2, three different pairs of rewards were shown, which mapped onto the same covariance structure, and the two contexts were interleaved within a run. An example ordering of runs is shown but this was balanced across participants.

(C) Participant view of exploring in one of the rooms during training. Floors of all rooms were purple in the scanner.

(D) Heatmaps of the average grid square occupancy per trial in each of the two contexts. Black arrows show the average transition vector from each grid square. Data are averaged across participants. Note that the average transition vectors were also well-matched when considering only those movement periods that were controlled by the participants (see [Figure S1D](#)).

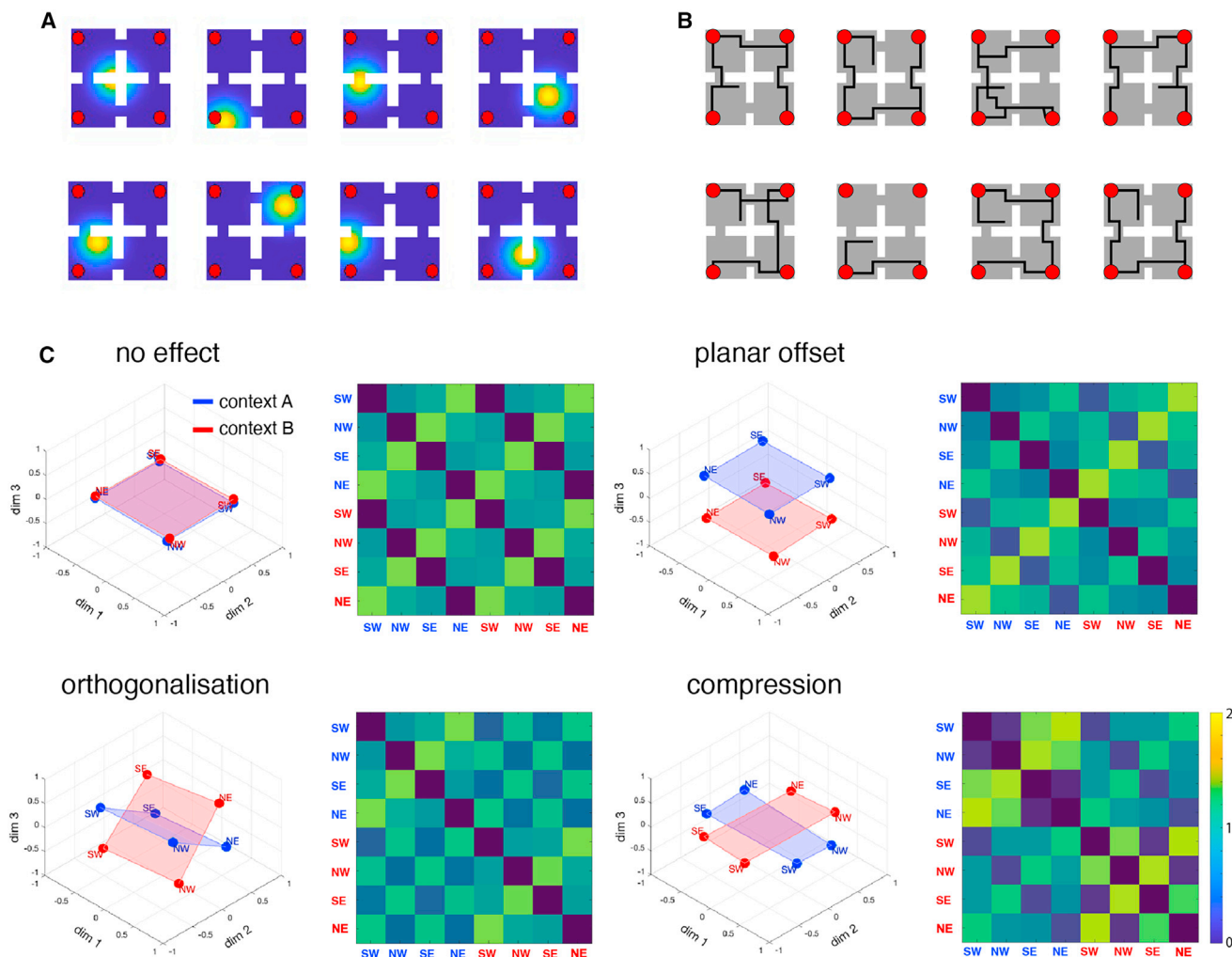
(E) Participant scores on each trial on days 1 (left) and 2 (right). With training, participants get faster at finding the rewards. On day 1, contexts were blocked across trials to facilitate learning, while on day 2 they were interleaved.

(F) Participants learn to preferentially search in rooms suggested by the reward structure, and this behavior generalizes to new sets of rewards associated with each context on day 2. In (D) and (E), data are shown smoothed across non-overlapping sets of 4 adjacent trials for visualization. In (E), we show only the room choices made by human participants and exclude those made by the agent. Error bars show standard error of the mean across participants. Color panels indicate different epochs with the same reward pairs.

behavioral trajectories, implying that any deviations from this prediction observed in BOLD data cannot be explained by imbalance in occupancy probabilities or transition frequencies.

However, we additionally equipped the model with three free parameters, corresponding to the hypotheses that context is encoded by (1) orthogonalization, (2) separation, or (3) compression of spatial representations. First, we allowed the fraction of cells  $\beta$  that remapped (i.e., changed their preferred spatial location) between contexts to vary. Global remapping ( $\beta = 1$ ) leads to full *orthogonalization* of the spatial map in each context; hence, it predicts that the two planes representing space in

each context should rotate to lie at  $90^\circ$  to one another ([Figure 2C](#), lower left panel). Second, we allowed a subset of cells to explicitly code for context along a dimension perpendicular to space and applied a freely varying gain factor  $\gamma$  to this neural activity, which creates a planar offset (or *separation*) between context representations ([Figure 2C](#), upper right panel). This model variant assumes that spatial and nonspatial variables are factorized within the neural code.<sup>10,42</sup> Finally, we assumed that place cells may jointly encode the avatar's current position and the prospective goal locations on each trial.<sup>38</sup> This was achieved with a final free parameter encoding the relative mixture weight  $\omega$ ,



**Figure 2. Computational modeling**

(A) Example simulated place fields in the model of the four rooms environment. Each panel is one neuron. White bars are walls. Red dots are boulders. The blue-yellow map shows neural tuning of a single neuron.

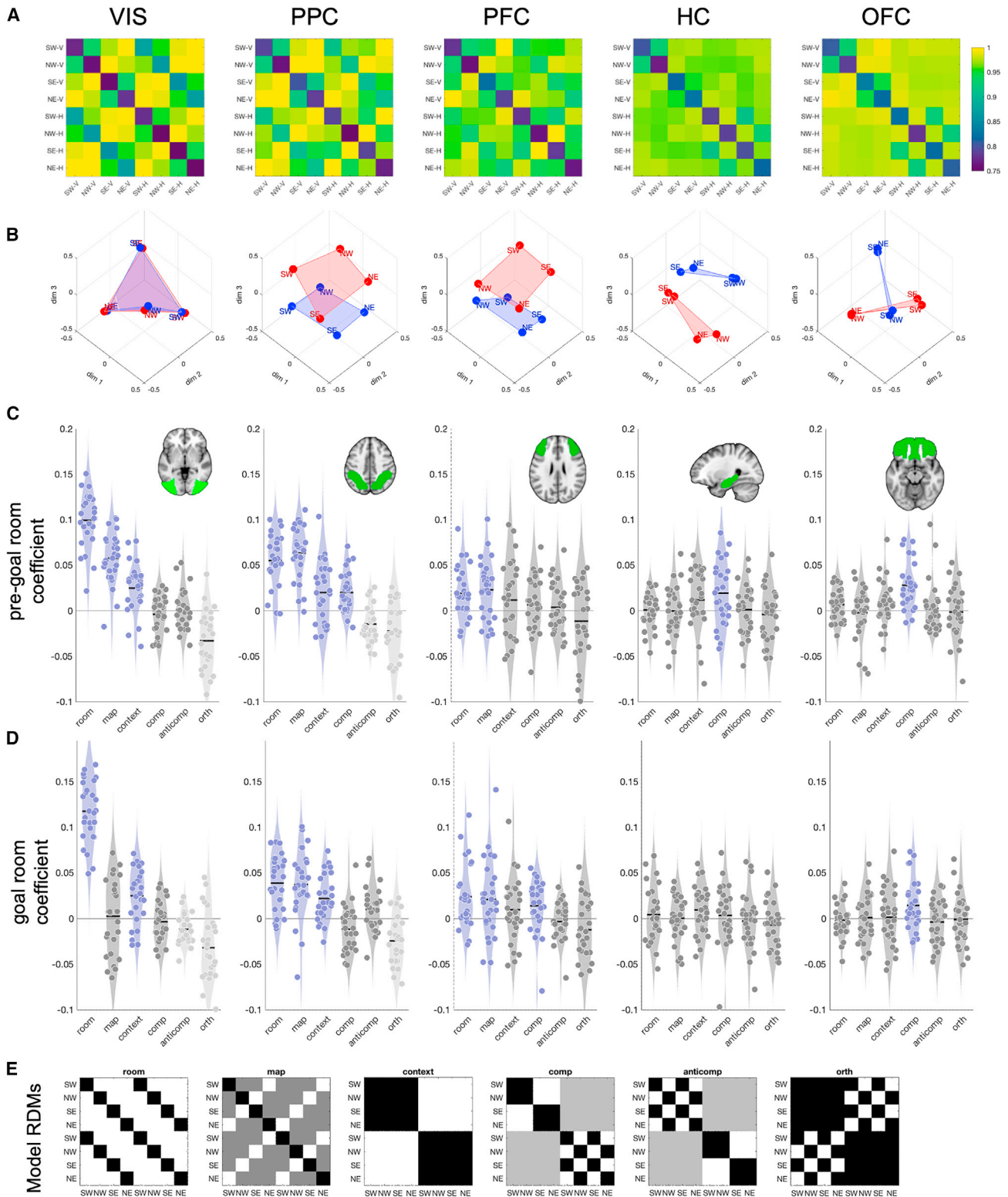
(B) Example trajectories made by a participant performing the task in the scanner (black lines). Red dots are boulders.

(C) Illustration of four representational hypotheses for different ways of separating two-dimensional (2D) information by context. On the left of each subpanel are MDS plots. Red and blue lines and shading indicate the different contexts. The dots denote the rooms (NE, NW, SE, SW) in each context. On the right of each panel is the corresponding RDM. Colors are in units of correlation distance. The data were generated under the following parameters: *no effect*,  $\beta = 0$ ,  $\gamma = 0.1$ ,  $\omega = 0$  (a small offset is introduced for ease of visualization); *planar separation*,  $\beta = 0$ ,  $\gamma = 0.5$ ,  $\omega = 0$ ; *orthogonalization*,  $\beta = 1$ ,  $\gamma = 0$ ,  $\omega = 0$ ; *compression*,  $\beta = 0$ ,  $\gamma = 0$ ,  $\omega = 0.9$ .

given to current and goal locations. Increasing the goal weight ( $\omega > 0$ ) leads to *compression*, whereby goal locations in a shared context are represented at lower distances than is warranted by their separation in physical space (i.e., the north and south rooms are closer together in the “vertical” context, and east and west in the “horizontal” context; Figure 2C, lower right panel). This occurs as the neural representations of current room and prospective goal are differentially mixed together in horizontal and vertical contexts. The goal weight parameter is designed to also allow the converse effect (anti-compression) when  $\omega < 0$ , which would be consistent with other recent observations.<sup>43</sup> Full details of the model are provided in the STAR Methods.

To test these hypotheses in humans, we estimated multivariate BOLD signals during navigation using a design matrix that modeled the presence of the avatar in each room (SW, NW, SE, and NE) and context (H and V) during the movement period yielding an  $8 \times 8$  RDM comparable with the model. All RDM analyses were conducted in cross-validation, comparing neural patterns between odd and even scanner runs. Goal approach is known to be a powerful modulator of BOLD signals,<sup>41,44,45</sup> and navigational choices are only made up until the point at which the goal room is entered; hence, we begin by focusing separately on the movement phases in which participants are approaching a room containing a goal (pre-goal room period)





**Figure 3. Regression analyses of neural geometries**

(A) Group average RDMs for each ROI. Each  $8 \times 8$  RDM is ordered (SW, NW, SE, NE) for first the vertical and then the horizontal context. Warmer colors indicate greater dissimilarity, and cooler colors greater similarity.

(legend continued on next page)

and where they are inside a room that contains a goal (goal-room period). We focus on anatomically defined regions of interest (ROIs) that have previously been implicated in goal-directed behavior, including the PFC, posterior parietal cortex (PPC), HC, and OFC, as well as a control ROI in the visual cortex (shown inset in [Figure 3C](#)). For each region, we plotted RDMs ([Figure 3A](#)) and visualized neural geometries in three dimensions, using MDS ([Figure 3B](#)). Finally, we complement this approach by presenting data from whole-brain searchlight analyses.

Visual inspection of the RDMs and MDS plots revealed that BOLD signals in these ROIs coded space and context in different ways (see [Videos S1](#), [S2](#), and [S3](#) for a clearer visualization in three dimensions [3D]). First, visual cortex represented each room with a distinct neural code, but the similarity structure was only weakly related to the overall spatial layout, and no effect of context was observed ([Figure 3B](#), far left panel). In PPC and PFC, rooms were encoded at the apices of two roughly parallel quadrilateral planes, thus with a geometry roughly matching the task environment (middle left and middle panels). Finally, in HC and OFC, spatial representations were compressed along the task-irrelevant axis so that “north” and “south” rooms are coded as adjacent in the vertical context, and “east” and “west” are represented as adjacent in the horizontal context. This compression effect was accompanied by a weak context-dependent separation, in which the two contexts were divided along another neural dimension running perpendicular (at 90°) to that encoding allocentric space ([Figure 3B](#), middle and far right panels).

To quantify these effects, we constructed model RDMs and regressed them against the data RDM in each ROI. We used 6 model RDMs in total, and these were included competitively in the regression. The first two model RDMs related to the structure of the environment. The (1) *room* model encoded each individual room with a unique code; the (2) *map* model encoded residual similarity structure that reflected the organization of the four rooms into a regular quadrilateral. The remaining RDMs encoded bases for (3) *separation* between contexts as predicted by  $\gamma > 0$ ; (4) the additional effects of both *compression* and (5) *anti-compression*, as predicted by  $\omega > 0$  and  $\omega < 0$ , respectively; and (6) the effect of *orthogonalization* as predicted by  $\beta > 0$ . Results are shown in [Figure 3C](#) for the pre-goal room period and in [Figure 3D](#) for the goal-room period (corresponding RDMs and MDS for the latter case are shown in [Figures S1A–S1C](#)). The relevant statistics (corrected for multiple comparisons) are reported in [Table 1](#); we use a false discovery rate correction for multiple comparison across independent ROIs. The model RDMs are also depicted in [Figure 3E](#).

### Encoding of spatial layout

We first considered how the spatial layout of the environment was coded in BOLD signals. Our model RDMs included predictors based on unstructured room identity (*room*) and residual

structure indicating the geometry of the environment (*map*). The first observation was that layout was coded most reliably in PPC and PFC. There was also a tendency for *room* to be more robustly encoded than *map* in visual cortex, especially during the goal-room period (where *map* accounted for no residual variance in visual cortex) but even during the pre-goal room period *room* explained numerically more variance than *map* in visual cortex. This can also be seen in the MDS plots in [Figure 3C](#), where the representations of the four rooms are roughly planar in PPC but in visual cortex are folded into a tetrahedron (on a 3D simplex) to accommodate the equal similarity between rooms. Thus, there appears to be a rough progression from a more unstructured, high-dimensional representation of the spatial layout (in visual cortex) to one in which rooms lie on quadrilateral planes in neural space, mirroring the layout of the four room environment (in PPC).

### Goal-based spatial compression

In HC and OFC, we did not observe a neural representation of the veridical spatial layout of the environment. Rather, we saw an effect of goal-based spatial compression, whereby rooms that were linked by virtue of shared goals were coded together. This is the pattern predicted by parameterizations of our model in which  $\omega > 0$ , i.e., where the agent’s location and goal are jointly coded in BOLD signals. This is clearly visible in the MDS plots, where (for both HC and OFC) the north and south rooms are more proximal in the vertical context, and east and west rooms are more proximal in the horizontal context ([Figure 3C](#)). Compression regressors were significant in both periods for OFC and the pre-goal room period for the HC, whereas effects of planar separation are weak or marginal in these ROIs ( $t$  values  $< 2$  in HC and OFC), and effects of orthogonalization are not significant. We did not observe differences in the strength of the compression effect in medial and lateral subportions of OFC, and the pre-goal room effect was independently significant in each subregion. As a control, we also verified that these effects did not occur prior to the first boulder being reached, at which point participants cannot know the optimal trajectory on that trial ([Figures S6A–S6C](#)).

### Model-free analyses

To complement these analyses based on regression models, which can be difficult to interpret especially when there is partial collinearity between predictors, we adopted an approach that involved averaging selected distances (vertex pairs) across data RDMs to ask targeted questions about the neural geometry. To achieve this, we constructed “score matrices” indicating which pairs of vertices were compared with each other; these are shown in [Figure 4A](#). The results largely mirrored those for regression analyses, with strong compression observed in HC and OFC during the pre-goal room period (see [Table S1](#)). This

(B) MDS plots (from the group average RDM) for each region. Blue dots are rooms in the vertical context and red in the horizontal context. For legibility, cardinally adjacent rooms within a context are linked by lines, which collectively form a quadrilateral when allocentric space is coded in just 2 dimensions.

(C) Violin plots showing coefficients for a competitive regression of model RDMs against each data RDM for the pre-goal room period. Each participant is an individual dot. Blue dots and shading (positive values) and light gray dots and shading (negative values) indicate  $p < 0.05$ .

(D) Same as (C) but for the goal-room period.

(E) Visualization of model RDMs used for these analyses; lighter colors indicate greater dissimilarity.

**Table 1. Statistics on regression coefficients**

| Pre-goal room period | Room    | Map     | Separation | Compress | Anti-compress | Orthogonalize |
|----------------------|---------|---------|------------|----------|---------------|---------------|
| VIS                  | 17.3*** | 10.7*** | 4.75***    | -1.04    | 0.31          | -5.44         |
| PPC                  | 9.34*** | 11.5*** | 3.06**     | 4.68***  | -3.95         | -3.61         |
| PFC                  | 3.50**  | 3.72*** | 1.50       | 1.13     | 0.72          | -1.44         |
| HC                   | 0.22    | -0.02   | 1.83       | 2.92**   | 0.24          | -0.75         |
| OFC                  | 2.04    | -0.42   | 1.94       | 5.70***  | 0.43          | -0.24         |
| Goal-room period     | Room    | Map     | Separation | Compress | Anti-compress | Orthogonalize |
| VIS                  | 18.3*** | 0.32    | 4.45***    | -0.80    | -3.69         | -4.46         |
| PPC                  | 6.77*** | 5.25*** | 4.61***    | -2.42    | 2.13          | -4.49         |
| PFC                  | 3.62*** | 2.60**  | 1.43       | 2.57**   | -0.81         | -1.67         |
| HC                   | 0.72    | 0.09    | 1.88       | 0.57     | -0.47         | -1.40         |
| OFC                  | -0.47   | 0.22    | 0.23       | 2.66**   | -0.84         | -0.20         |

t values for a test of each model RDM against zero for the data RDM from the pre-goal room period (upper part) and goal-room period (lower part). Each row is a brain region, and each column is a predictor. Asterisks: \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$  after false discovery rate (FDR) correction. According to a Shapiro-Wilks test, the OFC data were not normally distributed ( $p = 0.03$ ); so, we additionally conducted a non-parametric (sign) test against zero; the p value associated with this test was  $p < 0.001$ . VIS, visual cortex.

implies that the compression observed in the regression analysis was not an artifact of the other included predictors. We also used the place cell model to create model RDMs based on the best-fitting variants of models in which each of the three parameters  $\beta$ ,  $\gamma$ , and  $\omega$  were allowed to vary (or none of the three). This also confirmed that the neural data were best explained by a compression-based account in both HC and OFC (Figure S2).

### Correlations between behavior and brain activity

Next, we examined how across-cohort variation in compression scores for each brain region related to individual differences in behavior. One way to characterize individual participant performance is *transition bias*, which is the relative fraction of transitions made horizontally and vertically between rooms in the H and V contexts (a player that understands the structure should make proportionally more horizontal transitions in H context and vertical in the V context). An alternative measure is *first-choice accuracy*, which indexes whether participants' first transition reveals that they understand the correlation structure of the spatial goals in each context. For completeness, we correlated these behavioral measures with *compression*, *separation*, and *map* scores, although our main prediction was that compression would covary with performance in HC and OFC.

We observed that in the hippocampus, compression scores from the pre-goal room period positively predicted both transition bias ( $r = 0.45$ ,  $p = 0.019$ ) and first-choice accuracy ( $r = 0.43$ ,  $p = 0.025$ ). By contrast, in the OFC, compression score from the goal room period positively predicted first-choice accuracy in the goal room period ( $r = 0.41$ ,  $p < 0.032$ ). We plot the results of this correlation for HC and OFC in Figure 4D; results for other regions are shown in Figure S3.

### Inverted neural geometry across periods

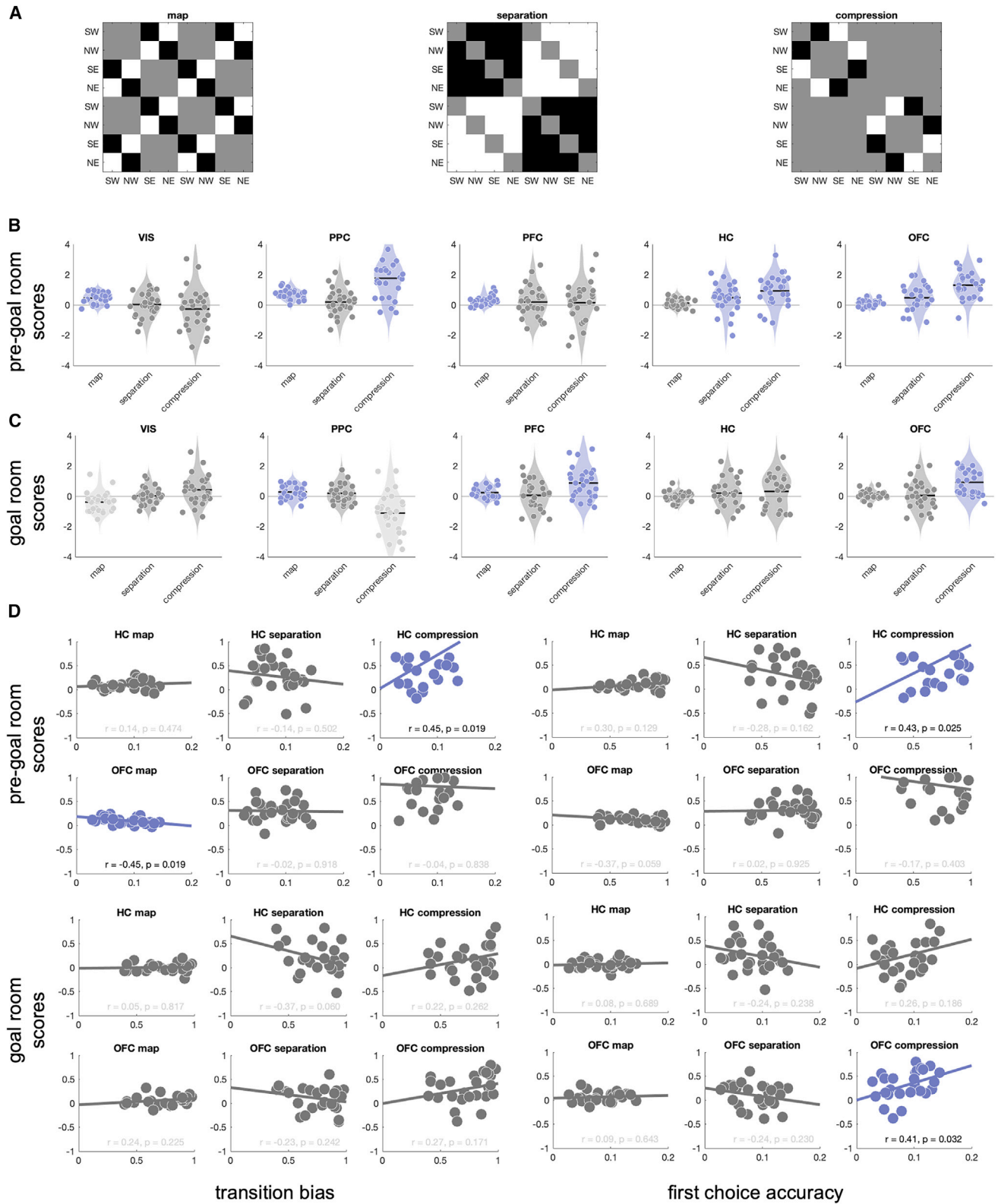
Having examined the geometries for the pre-goal room and goal-room periods, we next explored how they relate to one another. RDMs and corresponding MDS plots for the full period  $\times$  room  $\times$  context analysis are shown in Figures 5A and 5B. In the MDS plots, the pre-goal room period is now shown in cyan

(cue V) and orange (cue H), and the goal-room period in blue (cue V) and red (cue H). As can be seen, the brain powerfully encodes whether the agent is currently occupying a room with a goal, visible as the checker pattern in the RDMs and the resulting one-dimensional offset between periods that lies along a neural dimension perpendicular to that coding allocentric space. These results are confirmed by regressing model RDMs against the full  $16 \times 16$  data RDM (Figure S4); the effect of *period* was highly significant in each region (all t values  $> 13$ , all p values  $< 0.001$ ), with other effects mostly mirroring those described above (Table 2). Note that all analyses are conducted in cross-validation; hence, results are unlikely to be spuriously driven by temporal autocorrelation in BOLD signals. It is, however, consistent with previous reports that BOLD signals are powerfully modulated on the approach to a goal.<sup>41,44</sup> We can see that the effects of spatial layout (with or without compression) are thus represented in two parallel planar geometries, with a large offset coding whether the agent is currently occupying the goal room or is still navigating toward it. Interestingly, in the MDS plots for OFC and HC, the orientation of the planes for contexts H and V appears flipped between the two contexts, such that the coding of space and context is inverted when it is held in memory (during the pre-goal room period) and when it is being executed (during the goal-room period).

### Cosine similarity of neural vectors within and between periods

To quantify this latter effect, we computed the angle of the high-dimensional neural vector between each room/context and every other, both within periods (e.g., goal room to goal-room period) and across periods (e.g., goal room to pre-goal room period). Assuming a stylized model in which contexts were represented as compressed planes that were offset and inverted between periods (Figure 5C), we averaged angles for those edges that the model predicted to be parallel (e.g., common directions within a context; orange lines), orthogonal (e.g., perpendicular directions within a context; purple lines), and inverted (e.g., common directions within a context, but across periods;





(legend on next page)

cyan lines). Note that this analysis was conducted in the high-dimensional space of neural activity, not its compression to 3D in the MDS plots. In Figure 5D, we plotted the resulting angles for each ROI, which range from fully parallel (0) to fully orthogonal ( $\pi/2$ ) to inverted ( $\pi$ ). The results show that across regions, but especially for HC and OFC, there is a bias for edges within a context (e.g., NE-SE) to be parallel with edges denoting a common direction in space (e.g., NW-SW); that by contrast, those edges were more orthogonal to those denoting a perpendicular direction in space (e.g., NE-NW); and that edges that denoted a common direction across periods (e.g., NE-SE in pre-goal room period and NW-SW in goal-room period) had an even greater angular separation. We observed a main effect of the region ( $F_{3.6, 93.6} = 9.91, p < 0.001$ ) and a region  $\times$  pair type interaction ( $F_{5.9, 153.0} = 20.8, p < 0.001$ ), with the strongest separation of angle between edge pair types in the HC and OFC relative to other regions. These results thus confirm that the neural vectors respect the geometry of the environment within each period, but are inverted between periods, especially in HC and OFC

### Searchlight analyses

The foregoing analyses all rely on 5 ROIs that we chose *a priori*, given their previously described involvement in context-sensitive decision-making, navigation, and planning. However, to study these effects at the whole-brain level, we combined the score analysis with a whole-brain searchlight approach, allowing us to render map, separation, and compression effects onto a template brain. Results are consistent with our ROI analyses, and all regions described here contain searchlights, which reach significance at the whole-brain family-wise error-corrected [FWE] level. Figure 5F shows a visualization of the searchlight results for the pre-goal room period at a slightly more liberal statistical threshold ( $p < 0.0001$ , uncorrected) to facilitate the illustration of smaller clusters (see Figure S5 for FWE corrected whole-brain maps and further detail).

### Neural geometry of current and prospective locations

Finally, we asked how spatial goals were represented in each ROI, and how their representational geometry related to that of the current location in space. To ensure sufficient trial counts for this analysis, we collapsed over context, and modeled the BOLD data at the first level general linear model (GLM) with regressors coding for the currently occupied room and the location of the current goal, in a  $4 \times 4$  factorial design. This allowed us to construct model RDMs that encoded allocentric space as individual rooms or as a map (*room* and *map*, exactly as above) alongside new model RDMs that encoded current goal locations

as individual rooms or as a map (*goalroom* and *goalmap*; Figure 6A). Regression against the  $16 \times 16$  data RDM revealed that *goalmap* was significant in visual cortex, PPC and PFC, in addition to *map*. No effects were significant in HC or OFC, presumably because collapsing over orthogonal contexts removed the relevant subspaces in which rooms and goals are represented.

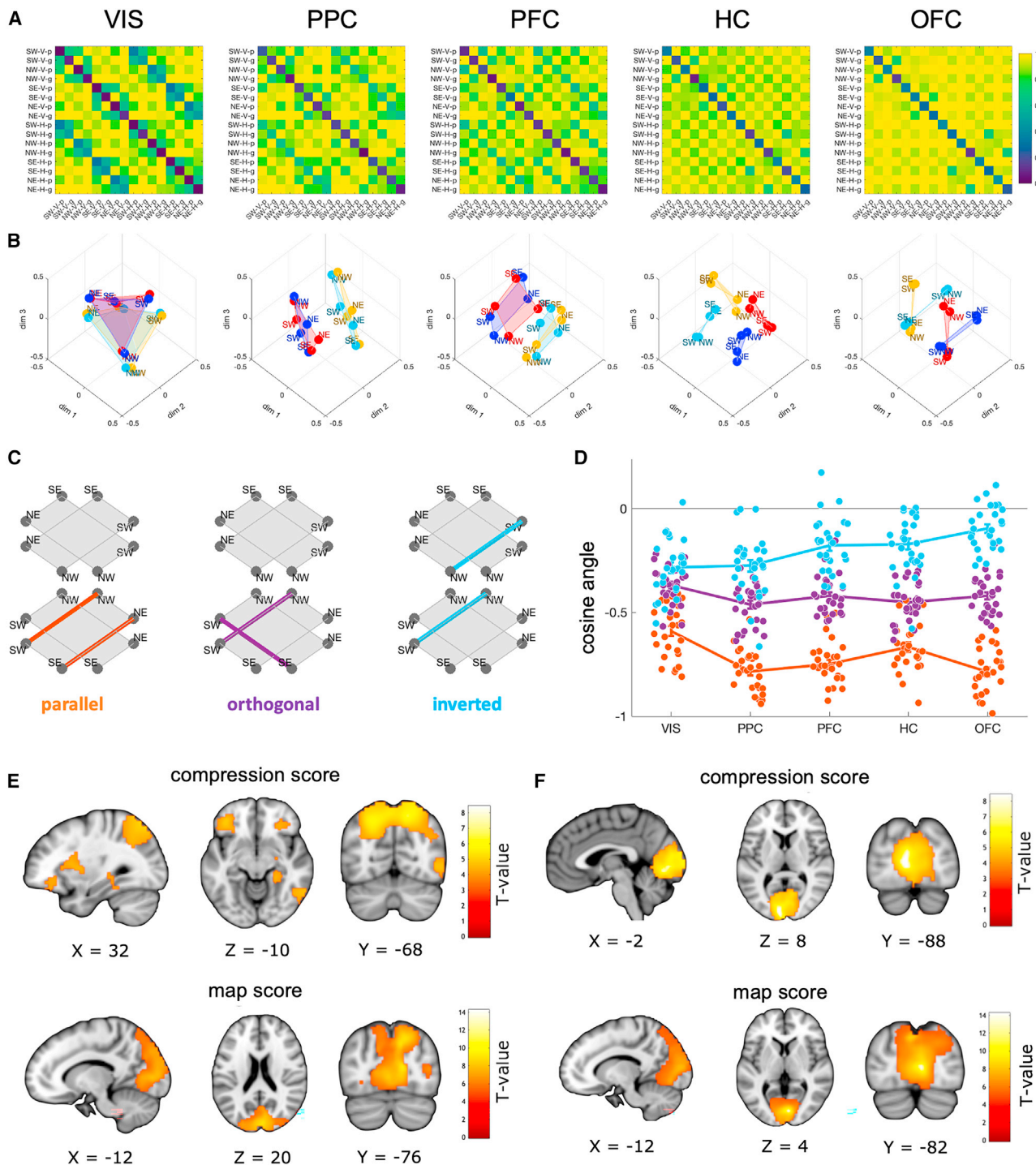
The full data from the regression analysis are shown in Figure 6B. RDMs for this analysis are shown in Figure 6C, and the MDS plots in Figure 6D. To increase legibility, we invert the plotting convention of the previous analysis and now plot different rooms (in allocentric space) in different colors (blue, SW; red, NW; cyan, SE; and orange, NE), and the labels on the plot now refer to goals (where the agent is headed). In PPC and PFC, goals are represented on 4 rough quadrilaterals, one within each room that the agent could occupy, although the representation of goals is smaller in area than the representations of room in allocentric space. There is thus a clear hierarchical representation, whereby a map of the current goal is represented within a map of the current location. In PPC, the quadrilateral is visibly elongated so that the goal-room condition is represented on a common plane separated from the non-goal-room conditions. The pattern in visual cortex is harder to discern. In PPC and PFC, thus, spatial goals are represented in a geometric format similar to physical space itself.

### DISCUSSION

Our major question was how context modulates the neural representation of allocentric space in human BOLD signals. We considered three major hypotheses. First, we asked whether context would lead to remapping, whereby population codes for space change their tuning preferences between contexts. This is implied by previous work in rodents showing that changes to the physical nature of the environment, or even changes in internal variables, can cause the spatial preferences of place cells to randomly remap.<sup>19–22,26</sup> One salient observation in the current report is that although context provoked representational changes in BOLD signals in hippocampus and neocortex, none of these changes resembled those expected if neural codes for space randomly remap, either partially or in full. The neural geometry implied by random remapping is that spatial representations become “orthogonal” or uncorrelated. By contrast, we observed that neural manifolds representing space were highly aligned across contexts in most brain regions. This resembles the “neural structure alignment” that has recently been reported to accompany decision tasks in both humans and monkeys,

### Figure 4. Score analyses and correlations with performance

- (A) Matrices used to compute scores. White entries are positive values (+1), black entries are negative values (−1), and gray entries are zeros (ignored). Each matrix was multiplied elementwise with the data RDM, and the resulting values summated to compute the corresponding score. The score matrices for *compression*, *separation*, and *map* are fully orthogonal.
- (B) Violin plots showing scores (*map*, *separation*, and *compression*) for the pre-goal room period in each region. Each dot is an individual participant. Blue dots and shading (positive values) and light gray dots and shading (negative values) indicate  $p < 0.05$ .
- (C) Same as (A) but for the goal-room period. Below, the (C) left: correlations between each score (see plot title) and transition bias for the HC (upper panels) and OFC (lower panels). Right: the same plots for transition bias. Blue dots denote significant ( $p < 0.05$ ) correlation.
- (D) Correlations between neural scores (*map*, *separation*, and *compression*) and behavioral measures (transition bias and first-choice accuracy) for HC and OFC in the pre-goal room period (upper panels) and goal-room period (lower panels). Each dot is a single participant, and the line is the best linear fit. Blue coloring is used to highlight significant correlations ( $p < 0.05$ ).



**Figure 5. Inverted neural geometries across periods and searchlight analyses**

(A) Group average RDMs for the full period  $\times$  room  $\times$  context analysis in each region. Each RDM comprises the nested variables period (goal room, pre-goal room), room (SW, NW, SE, NE), and context (vertical, horizontal). Warmer colors signal greater dissimilarity.

(B) MDS plots from the corresponding group average RDM. Colors denote period/context combinations: dark blue, vertical, pre-goal room; red, horizontal, pre-goal room; cyan, vertical, goal room; orange, horizontal, goal room.

(C) Stylized model of the MDS plots in (B), to illustrate those edges predicted to be parallel (left panel, orange lines), orthogonal (middle panel, purple lines), and inverted (right panel, cyan lines).

(D) Cosine angle between neural vectors for all predicted parallel, orthogonal, and inverted edges, rendered onto a single plot. Dots are individual participants, and the line shows the average for each region.

(legend continued on next page)

**Table 2. Statistics on regression coefficients from a regression modeling both pre-goal and goal-room periods**

|     | Room     | Map     | Period   | Separation | Compression (goal room) | Compression (pre-goal room) | Orthogonalize |
|-----|----------|---------|----------|------------|-------------------------|-----------------------------|---------------|
| VIS | 18.35*** | 2.84**  | 20.6***  | 0.60       | 2.28                    | −1.00                       | −4.50         |
| PPC | 9.52***  | 8.76*** | 20.7***  | 0.63       | −4.33                   | 7.06***                     | −1.39         |
| PFC | 3.90***  | 5.08*** | 16.7***  | 1.04       | 4.33                    | 0.58                        | −0.21         |
| HC  | −0.30    | −1.47   | 13.0***  | 2.91**     | 1.53                    | 4.82***                     | 1.90          |
| OFC | −1.91    | −7.46   | 14.45*** | 1.76       | 5.97***                 | 7.85***                     | 3.98***       |

t values for a test of each model RDM coefficient against zero for the data RDM from the goal-room period. Each row is a brain region, and each column is a predictor. Asterisks: \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ .

whereby contexts sharing common structure are represented with parallel neural geometries, potentially because this allows a decoder trained in one context to be generalized to the other.<sup>46–51</sup>

The second hypothesis we considered was that context is represented as an independent, nonspatial dimension in neural state space. This is implied by recent findings emphasizing that multiple task-relevant variables, such as location and evidence for progress toward a goal, are multiplexed in neurons with spatial selectivity, e.g., in the rodent hippocampal CA1 area.<sup>42</sup> Although participants were navigating toward a room containing a goal, we did see evidence for a nonspatial representation of context in both hippocampus and OFC, as evidenced by a reliable offset or “separation” between contexts in the neural manifolds for space. However, this effect was less prominent and statistically weaker than the other effects reported here, did not survive whole-brain correction, and did not persist in either HC or OFC once participants entered the goal room.

Instead, in our study, the most salient way that context influenced neural coding was by compressing spatial codes so that prospective locations signaled by the context lay closer together in neural state space. Compression effects were most prominent in the OFC and hippocampus, where the rooms lying along the relevant dimensions were coded with a highly correlated neural code. Using a spatial encoding model, we show that it is possible to elicit compression of this sort via the simple assumption that both current and prospective locations are encoded jointly in population vectors for allocentric space. This occurs because jointly encoding prospective locations (the two spatial goals) that lie on a common axis leads to correlated neural signals along this axis, which in turn are visible in compressions of the neural geometry for space.

The compression in HC and OFC was sufficiently prominent that in our context-dependent navigation task, neither region naively reliably encoded the full spatial layout, as might be expected from a pure place code. This might seem curious, given that the hippocampus encodes a spatial representation of the environment in rodents. However, there are a number of possible explanations for this. First, it seems likely that compression is incurred by the need to keep representations of possible spatial

plans separate: to ensure that horizontal and vertical goals are not confused. In which case, it is possible that compression does not occur in standard navigation paradigms where goals are not flexibly cued from trial to trial and that there the HC and OFC maps resemble more closely those observed in PPC. Second, it is possible that there are variations in the extent to which current and goal locations are decodable from human BOLD signals relative to neuronal recordings in rodents. Indeed, prospective information seems to be a prominent component of human BOLD responses in a variety of settings,<sup>9,17,34,52,53</sup> and the place code seen ubiquitously in rodents is much less prominent in monkeys<sup>54</sup> and humans.<sup>55</sup> It is also possible that this is due to differences in recording methods; there is considerable debate about how to jointly understand effects recorded at the micro-, meso-, and macro-scopic levels during spatial navigation.<sup>56</sup>

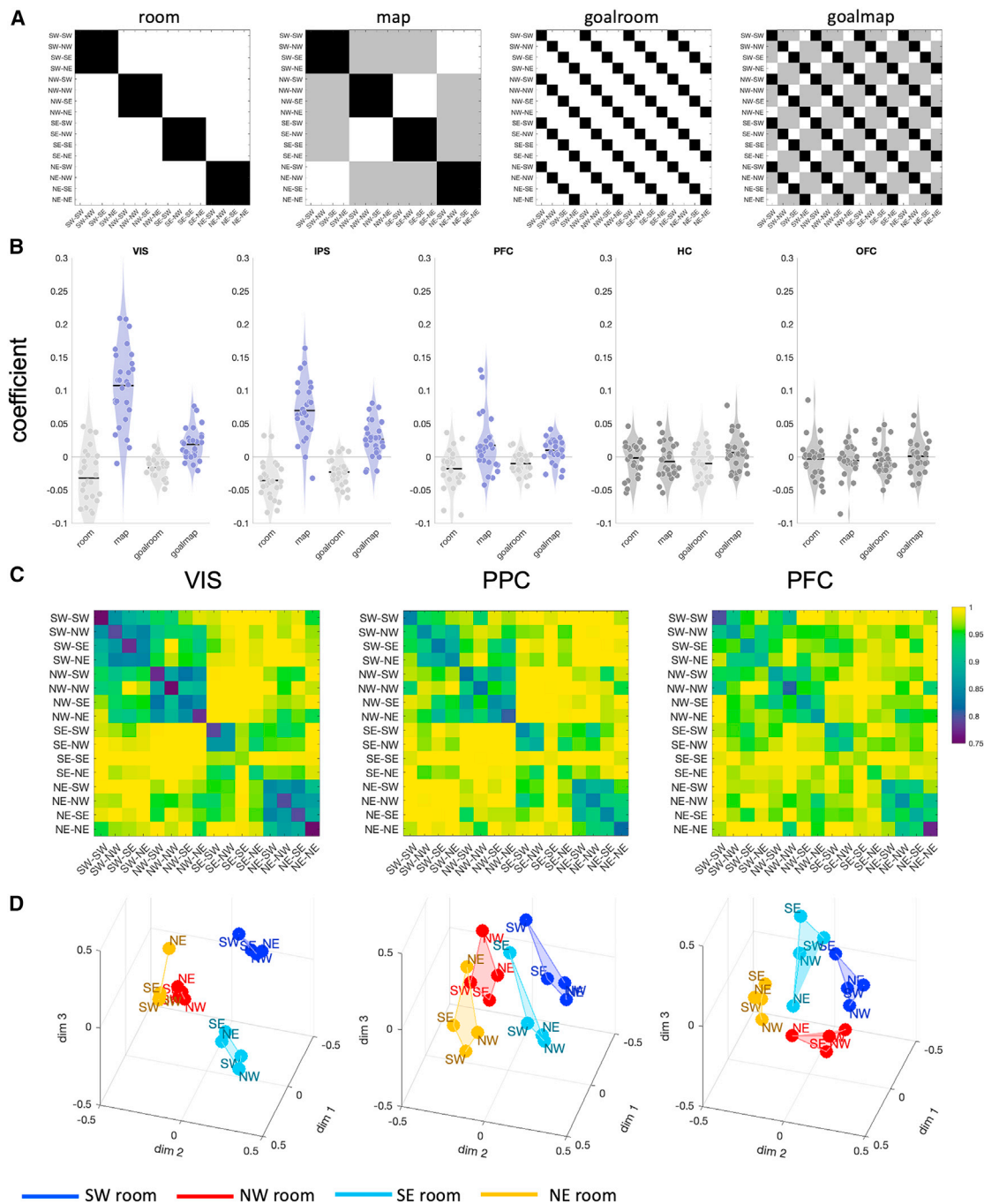
Our results are thus consistent with the finding that both hippocampus<sup>9,13,34,37,38</sup> and orbitofrontal cortex<sup>40,57</sup> explicitly code for future goal locations. Our model suggests that the representation of space in the BOLD signal in hippocampus and OFC can be explained by the simple principle that current and prospective (goal) locations are encoded in temporal proximity, but our recording methods do not have the resolution in space or time to detail exactly how that might occur. For example, prospective locations or goals may be represented through dedicated cell types<sup>38</sup> or may be evoked during forward or backward simulation occurring via replay mechanisms,<sup>58,59</sup> which has also been observed in humans.<sup>60</sup> More generally, our results are consistent with the view that the OFC (and to a lesser extent HC) represents the “task space,” that is, it encodes states in a format that is optimized for reward-guided action and planning.<sup>61,62</sup> This could explain the context-dependent compression of vertical/horizontal rooms, which represent the reward-relevant axes of our task.

We observed another curious effect by which the neural geometry of the environment was “flipped” between periods in which (1) navigation was ongoing and (2) where the goal room had been reached. It is not clear to us what purpose is served by this aspect of the geometry, which was most prominent in HC and OFC. However, it is reminiscent of recent reports that

(E) Searchlight analyses: whole-brain effects of *compression score* and *map score* for the pre-goal room period, rendered onto a template brain at a threshold of  $p < 0.0001$  uncorrected. All regions shown contain voxels significant at  $p < 0.05$  after family-wise error correction (whole-brain images thresholded with family-wise error correction at  $p < 0.05$  are shown in the [supplemental information](#)).

(F) Same as (E) but for goal room period.





**Figure 6. Neural geometry of current and prospective locations**

(A) Model RDMs used to test the neural geometry of room and goal representations. RDMs are constructed from nested goal (SW, NW, SE, NE) and room (SW, NW, SE, NE) variables.

(B) Violin plots showing parameter estimates for each model RDM regressed competitively against the data RDM. Each dot is an individual participant. Blue dots and shading (positive values) and light gray dots and shading (negative values) indicate  $p < 0.05$ .

(C) Group average data RDMs for visual cortex, PPC and PFC (HC and OFC showed no significant effects).

(D) MDS plots constructed from corresponding group average RDM for each region. Colors denote room (blue, SW; red, NW; cyan, SE; orange, NE). Rooms are organized into an approximate quadrilateral, and goals (within each room) are similarly arranged approximately quadrilaterally.

memory traces are rotated in neural state space to prevent them from interfering with perceptual information.<sup>63</sup> Other reports have emphasized a related effect, whereby retrieval induces spatial memories to be mutually repulsed and become more distinguishable.<sup>43</sup> It seems possible that retrieval-based repulsion (between periods) and context-based compression of spatial memories can co-exist; this may be an interesting avenue for future research.

We also examined how the agent's current location and the location of the navigational goal were jointly represented. Remarkably, we observed that the representation of agent location and goal location is nested, especially in PPC. Here, we observed a prominent quadrilateral representation of the occupied room, but nested within each room representation was another quadrilateral representation of the navigational goal. Consistent with the strong effect of period, this representation is distorted (at least in PPC), so that the goal corresponding to the current room is represented distinctly from all other goals. This hierarchical representation of goals and space was not observed in HC or OFC in our study, presumably because averaging over contexts removes the subspace in which location is represented.

How context biases the encoding of sensory signals has been extensively studied in tasks that require a single action to be taken to elicit an outcome, such as visual categorization. In these tasks, different computational mechanisms have been proposed for preventing interference between different tasks (or goals) that are required in different contexts. For example, when the task requires monkeys to classify stimuli into common groups, single neurons, or populations in PFC code for stimuli associated with a given class,<sup>6,64–66</sup> echoing the compression of target information reported here. Other reports, however, argue that during categorization, neural signals coding for different groups are offset by a one-dimensional signal, giving rise to a neural separation similar to that tested here.<sup>67,68</sup> Where there are explicit contextual cues signaling the task, context-irrelevant information can be compressed in BOLD signals<sup>7</sup> and is often coded along a perpendicular dimension in neural state space, e.g., to avoid catastrophic interference.<sup>6,7</sup> Thus, orthogonalization, separation, and compression are candidate mechanisms for mediating the contextual modulation of sensory codes in both instantaneous and sequential decision tasks.

## STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- **KEY RESOURCES TABLE**
- **RESOURCE AVAILABILITY**
  - Lead contact
  - Materials availability
  - Data and code availability
- **METHOD DETAILS**
  - Participants
  - Design, task and procedures
  - Behavioural analysis
  - Computational model

- fMRI data collection and pre-processing
- fMRI analysis: first level GLMs and ROI definition
- fMRI analysis: neural geometry

## SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.neuron.2023.08.021>.

## ACKNOWLEDGMENTS

This work was supported by a Sir Henry Wellcome Fellowship (award 210849/Z/18/Z) to P.S.M.-K., by funding from the European Research Council (ERC Consolidator awards 725937 to C.S. and 820213 to G.P. and ERC Starting grant 852669 to N.W.S.), by Special Grant Agreement No. 945539 (Human Brain Project SGA) to C.S., H.S., and G.P., and by funding from the Max Planck Society (Max Planck Research Group grant M.TN.A.BILD0004) and from the Federal Government of Germany and the Laender under the Excellence Strategy to N.W.S. The authors would like to thank Kate Jeffery, Eleonore Duvelle, and Roddy Grieves for helpful comments on this project.

## AUTHOR CONTRIBUTIONS

C.S., G.P., and H.J.S. originally conceived the study. H.S. and S.C. performed the experiments (under the supervision of C.S. and N.W.S.). P.S.M.-K., H.S., and C.S. analyzed the data, performed simulations and visualizations, and wrote the original draft of the paper. G.P., H.J.S., and N.W.S. reviewed and edited the original draft of the paper.

## DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: January 13, 2023

Revised: May 10, 2023

Accepted: August 18, 2023

Published: September 18, 2023

## REFERENCES

1. Miller, E.K., and Cohen, J.D. (2001). An integrative theory of prefrontal cortex function. *Annu. Rev. Neurosci.* 24, 167–202. <https://doi.org/10.1146/annurev.neuro.24.1.167>.
2. Siegel, M., Buschman, T.J., and Miller, E.K. (2015). Cortical information flow during flexible sensorimotor decisions. *Science* 348, 1352–1355. <https://doi.org/10.1126/science.aab0551>.
3. Ester, E.F., Sprague, T.C., and Serences, J.T. (2020). Categorical biases in human occipitoparietal cortex. *J. Neurosci.* 40, 917–931. <https://doi.org/10.1523/JNEUROSCI.2700-19.2019>.
4. Mante, V., Sussillo, D., Shenoy, K.V., and Newsome, W.T. (2013). Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature* 503, 78–84. <https://doi.org/10.1038/nature12742>.
5. Brincat, S.L., Siegel, M., von Nicolai, C., and Miller, E.K. (2018). Gradual progression from sensory to task-related processing in cerebral cortex. *Proc. Natl. Acad. Sci. USA* 115, E7202–E7211. <https://doi.org/10.1073/pnas.1717075115>.
6. Roy, J.E., Riesenhuber, M., Poggio, T., and Miller, E.K. (2010). Prefrontal cortex activity during flexible categorization. *J. Neurosci.* 30, 8519–8528. <https://doi.org/10.1523/JNEUROSCI.4837-09.2010>.
7. Flesch, T., Juechems, K., Dumbalska, T., Saxe, A., and Summerfield, C. (2022). Orthogonal representations for robust context-dependent task performance in brains and neural networks. *Neuron* 110, 1258–1270.e11. <https://doi.org/10.1016/j.neuron.2022.01.005>.
8. Takagi, Y., Hunt, L.T., Woolrich, M.W., Behrens, T.E., and Klein-Flügge, M.C. (2021). Adapting non-invasive human recordings along multiple

- task-axes shows unfolding of spontaneous and over-trained choice. *eLife* 10, e60988. <https://doi.org/10.7554/eLife.60988>.
9. Nyberg, N., Duvelle, É., Barry, C., and Spiers, H.J. (2022). Spatial goal coding in the hippocampal formation. *Neuron* 110, 394–422. <https://doi.org/10.1016/j.neuron.2021.12.012>.
  10. McKenzie, S., Frank, A.J., Kinsky, N.R., Porter, B., Rivière, P.D., and Eichenbaum, H. (2014). Hippocampal representation of related and opposing memories develop within distinct, hierarchically organized neural schemas. *Neuron* 83, 202–215. <https://doi.org/10.1016/j.neuron.2014.05.019>.
  11. Spiers, H.J., and Barry, C. (2015). Neural systems supporting navigation. *Curr. Opin. Behav. Sci.* 1, 47–55. <https://doi.org/10.1016/j.cobeha.2014.08.005>.
  12. O'Keefe, J., and Dostrovsky, J. (1971). The hippocampus as a spatial map. Preliminary evidence from unit activity in the freely-moving rat. *Brain Res.* 34, 171–175. [https://doi.org/10.1016/0006-8993\(71\)90358-1](https://doi.org/10.1016/0006-8993(71)90358-1).
  13. Ekstrom, A.D., Kahana, M.J., Caplan, J.B., Fields, T.A., Isham, E.A., Newman, E.L., and Fried, I. (2003). Cellular networks underlying human spatial navigation. *Nature* 425, 184–188. <https://doi.org/10.1038/nature01964>.
  14. Kunz, L., Staresina, B.P., Reinacher, P.C., Brandt, A., Guth, T.A., Schulze-Bonhage, A., and Jacobs, J. (2022). Ripple-locked coactivity of stimulus-specific neurons supports human associative memory. Preprint at bioRxiv. <https://doi.org/10.1101/2022.10.17.512635>.
  15. Hassabis, D., Chu, C., Rees, G., Weiskopf, N., Molyneux, P.D., and Maguire, E.A. (2009). Decoding neuronal ensembles in the human hippocampus. *Curr. Biol.* 19, 546–554. <https://doi.org/10.1016/j.cub.2009.02.033>.
  16. Kim, M., Jeffery, K.J., and Maguire, E.A. (2017). Multivoxel pattern analysis reveals 3D place information in the human hippocampus. *J. Neurosci.* 37, 4270–4279. <https://doi.org/10.1523/JNEUROSCI.2703-16.2017>.
  17. Rodriguez, P.F. (2010). Neural decoding of goal locations in spatial navigation in humans with fMRI. *Hum. Brain Mapp.* 31, 391–397. <https://doi.org/10.1002/hbm.20873>.
  18. Sulpizio, V., Committeri, G., and Galati, G. (2014). Distributed cognitive maps reflecting real distances between places and views in the human brain. *Front. Hum. Neurosci.* 8, 716. <https://doi.org/10.3389/fnhum.2014.00716>.
  19. Anderson, M.I., and Jeffery, K.J. (2003). Heterogeneous modulation of place cell firing by changes in context. *J. Neurosci.* 23, 8827–8835. <https://doi.org/10.1523/JNEUROSCI.23-26-08827.2003>.
  20. Leutgeb, S., Leutgeb, J.K., Barnes, C.A., Moser, E.I., McNaughton, B.L., and Moser, M.B. (2005). Independent codes for spatial and episodic memory in hippocampal neuronal ensembles. *Science* 309, 619–623. <https://doi.org/10.1126/science.1114037>.
  21. Sanders, H., Wilson, M.A., and Gershman, S.J. (2020). Hippocampal remapping as hidden state inference. *eLife* 9, e51140. <https://doi.org/10.7554/eLife.51140>.
  22. Markus, E., Qin, Y., Leonard, B., Skaggs, W., McNaughton, B., and Barnes, C. (1995). Interactions between location and task affect the spatial and directional firing of hippocampal neurons. *J. Neurosci.* 15, 7079–7094. <https://doi.org/10.1523/JNEUROSCI.15-11-07079.1995>.
  23. Dudchenko, P.A., and Wood, E.R. (2014). Splitter cells: hippocampal place cells whose firing is modulated by where the animal is going or where it has been. In *Space, Time and Memory in the Hippocampal Formation*, D. Derdikman and J.J. Knierim, eds. (Springer), pp. 253–272. [https://doi.org/10.1007/978-3-7091-1292-2\\_10](https://doi.org/10.1007/978-3-7091-1292-2_10).
  24. Ainge, J.A., Tamosiunaite, M., Woergoetter, F., and Dudchenko, P.A. (2007). Hippocampal CA1 place cells encode intended destination on a maze with multiple choice points. *J. Neurosci.* 27, 9769–9779. <https://doi.org/10.1523/JNEUROSCI.2011-07.2007>.
  25. Wood, E.R., Dudchenko, P.A., Robitsek, R.J., and Eichenbaum, H. (2000). Hippocampal neurons encode information about different types of memory episodes occurring in the same location. *Neuron* 27, 623–633. [https://doi.org/10.1016/S0896-6273\(00\)00071-4](https://doi.org/10.1016/S0896-6273(00)00071-4).
  26. Rennó-Costa, C., Lisman, J.E., and Verschure, P.F.M.J. (2010). The mechanism of rate remapping in the dentate gyrus. *Neuron* 68, 1051–1058. <https://doi.org/10.1016/j.neuron.2010.11.024>.
  27. Dupret, D., O'Neill, J., Pleydell-Bouverie, B., and Csicsvari, J. (2010). The reorganization and reactivation of hippocampal maps predict spatial memory performance. *Nat. Neurosci.* 13, 995–1002. <https://doi.org/10.1038/nn.2599>.
  28. Hollup, S.A., Molden, S., Donnett, J.G., Moser, M.B., and Moser, E.I. (2001). Accumulation of hippocampal place fields at the goal location in an annular watermaze task. *J. Neurosci.* 21, 1635–1644. <https://doi.org/10.1523/JNEUROSCI.21-05-01635.2001>.
  29. Grienberger, C., and Magee, J.C. (2022). Entorhinal cortex directs learning-related changes in CA1 representations. *Nature* 611, 554–562. <https://doi.org/10.1038/s41586-022-05378-6>.
  30. Hok, V., Lenck-Santini, P.P., Roux, S., Save, E., Muller, R.U., and Poucet, B. (2007). Goal-related activity in hippocampal place cells. *J. Neurosci.* 27, 472–482. <https://doi.org/10.1523/JNEUROSCI.2864-06.2007>.
  31. Ito, H.T., Zhang, S.J., Witter, M.P., Moser, E.I., and Moser, M.B. (2015). A prefrontal-thalamo-hippocampal circuit for goal-directed spatial navigation. *Nature* 522, 50–55. <https://doi.org/10.1038/nature14396>.
  32. Stachenfeld, K.L., Botvinick, M.M., and Gershman, S.J. (2017). The hippocampus as a predictive map. *Nat. Neurosci.* 20, 1643–1653. <https://doi.org/10.1038/nn.4650>.
  33. Frank, L.M., Brown, E.N., and Wilson, M. (2000). Trajectory encoding in the hippocampus and entorhinal cortex. *Neuron* 27, 169–178. [https://doi.org/10.1016/S0896-6273\(00\)00018-0](https://doi.org/10.1016/S0896-6273(00)00018-0).
  34. Brown, T.I., Carr, V.A., LaRocque, K.F., Favila, S.E., Gordon, A.M., Bowles, B., Bailenson, J.N., and Wagner, A.D. (2016). Prospective representation of navigational goals in the human hippocampus. *Science* 352, 1323–1326. <https://doi.org/10.1126/science.aaf0784>.
  35. Simon, D.A., and Daw, N.D. (2011). Neural correlates of forward planning in a spatial decision task in humans. *J. Neurosci.* 31, 5526–5539. <https://doi.org/10.1523/JNEUROSCI.4647-10.2011>.
  36. Kunz, L., Wang, L., Lachner-Piza, D., Zhang, H., Brandt, A., Dümpelmann, M., Reinacher, P.C., Coenen, V.A., Chen, D., Wang, W.X., et al. (2019). Hippocampal theta phases organize the reactivation of large-scale electrophysiological representations during goal-directed navigation. *Sci. Adv.* 5, eaav8192. <https://doi.org/10.1126/sciadv.aav8192>.
  37. Poucet, B., and Hok, V. (2017). Remembering goal locations. *Curr. Opin. Behav. Sci.* 17, 51–56. <https://doi.org/10.1016/j.cobeha.2017.06.003>.
  38. Gauthier, J.L., and Tank, D.W. (2018). A dedicated population for reward coding in the hippocampus. *Neuron* 99, 179–193.e7. <https://doi.org/10.1016/j.neuron.2018.06.008>.
  39. Burgess, N., and O'Keefe, J. (1996). Neuronal computations underlying the firing of place cells and their role in navigation. *Hippocampus* 6, 749–762. [https://doi.org/10.1002/\(SICI\)1098-1063\(1996\)6:6<749::AID-HIPO16>3.0.CO;2-0](https://doi.org/10.1002/(SICI)1098-1063(1996)6:6<749::AID-HIPO16>3.0.CO;2-0).
  40. Basu, R., Gebauer, R., Herfurth, T., Kolb, S., Golipour, Z., Tchumatchenko, T., and Ito, H.T. (2021). The orbitofrontal cortex maps future navigational goals. *Nature* 599, 449–452. <https://doi.org/10.1038/s41586-021-04042-9>.
  41. Howard, L.R., Javadi, A.H., Yu, Y., Mill, R.D., Morrison, L.C., Knight, R., Loftus, M.M., Staskute, L., and Spiers, H.J. (2014). The hippocampus and entorhinal cortex encode the path and Euclidean distances to goals during navigation. *Curr. Biol.* 24, 1331–1340. <https://doi.org/10.1016/j.cub.2014.05.001>.
  42. Nieh, E.H., Schottdorf, M., Freeman, N.W., Low, R.J., Lewallen, S., Koay, S.A., Pinto, L., Gauthier, J.L., Brody, C.D., and Tank, D.W. (2021). Geometry of abstract learned knowledge in the hippocampus. *Nature* 595, 80–84. <https://doi.org/10.1038/s41586-021-03652-7>.

43. Chanales, A.J.H., Oza, A., Favila, S.E., and Kuhl, B.A. (2017). Overlap among spatial memories triggers repulsion of hippocampal representations. *Curr. Biol.* 27, 2307–2317.e5. <https://doi.org/10.1016/j.cub.2017.06.057>.
44. Balaguer, J., Spiers, H.J., Hassabis, D., and Summerfield, C. (2016). Neural mechanisms of hierarchical planning in a virtual subway network. *Neuron* 90, 893–903. <https://doi.org/10.1016/j.neuron.2016.03.037>.
45. Bierbrauer, A., Kunz, L., Gomes, C.A., Luhmann, M., Deuker, L., Getzmann, S., Wascher, E., Gajewski, P.D., Hengstler, J.G., Fernandez-Alvarez, M., et al. (2020). Unmasking selective path integration deficits in Alzheimer's disease risk carriers. *Sci. Adv.* 6, eaba1394. <https://doi.org/10.1126/sciadv.aba1394>.
46. Bernardi, S., Benna, M.K., Rigotti, M., Munuera, J., Fusi, S., and Salzman, C.D. (2020). The geometry of abstraction in the hippocampus and prefrontal cortex. *Cell* 183, 954–967.e21. <https://doi.org/10.1016/j.cell.2020.09.031>.
47. Luyckx, F., Nili, H., Spitzer, B., and Summerfield, C. (2019). Neural structure mapping in human probabilistic reward learning. *eLife* 8, e42816. <https://doi.org/10.7554/eLife.42816>.
48. Sheahan, H., Luyckx, F., Nelli, S., Teupe, C., and Summerfield, C. (2021). Neural state space alignment for magnitude generalization in humans and recurrent networks. *Neuron* 109, 1214–1226.e8. <https://doi.org/10.1016/j.neuron.2021.02.004>.
49. Morton, N.W., Schlichting, M.L., and Preston, A.R. (2020). Representations of common event structure in medial temporal lobe and frontoparietal cortex support efficient inference. *Proc. Natl. Acad. Sci. USA* 117, 29338–29345. <https://doi.org/10.1073/pnas.1912338117>.
50. Nelli, S., Braun, L., Dumbalska, T., Saxe, A., and Summerfield, C. (2023). Neural knowledge assembly in humans and neural networks. *Neuron* 111, 1504–1516.e9. <https://doi.org/10.1016/j.neuron.2023.02.014>.
51. Ito, T., Klinger, T., Schultz, D.H., Murray, J.D., Cole, M.W., and Rigotti, M. (2022). Compositional generalization through abstract representations in human and artificial neural networks. *Adv. Neural Inf. Process. Syst.* 35, 32225–32239. <https://doi.org/10.48550/ARXIV.2209.07431>.
52. Schacter, D.L., and Addis, D.R. (2007). The cognitive neuroscience of constructive memory: remembering the past and imagining the future. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 362, 773–786. <https://doi.org/10.1098/rstb.2007.2087>.
53. Chadwick, M.J., Jolly, A.E., Amos, D.P., Hassabis, D., and Spiers, H.J. (2015). A goal direction signal in the human entorhinal/subicular region. *Curr. Biol.* 25, 87–92. <https://doi.org/10.1016/j.cub.2014.11.001>.
54. Killian, N.J., Jutras, M.J., and Buffalo, E.A. (2012). A map of visual space in the primate entorhinal cortex. *Nature* 491, 761–764. <https://doi.org/10.1038/nature11587>.
55. Nolan, C.R., Vromen, J.M.G., Cheung, A., and Baumann, O. (2018). Evidence against the detectability of a hippocampal place code using functional magnetic resonance imaging. *eNeuro* 5, ENEURO.0177-18.2018. <https://doi.org/10.1523/ENEURO.0177-18.2018>.
56. Kunz, L., Maidenbaum, S., Chen, D., Wang, L., Jacobs, J., and Axmacher, N. (2019). Mesoscopic neural representations in spatial navigation. *Trends Cogn. Sci.* 23, 615–630. <https://doi.org/10.1016/j.tics.2019.04.011>.
57. Feierstein, C.E., Quirk, M.C., Uchida, N., Sosulski, D.L., and Mainen, Z.F. (2006). Representation of spatial goals in rat orbitofrontal cortex. *Neuron* 51, 495–507. <https://doi.org/10.1016/j.neuron.2006.06.032>.
58. Pfeiffer, B.E., and Foster, D.J. (2013). Hippocampal place-cell sequences depict future paths to remembered goals. *Nature* 497, 74–79. <https://doi.org/10.1038/nature12112>.
59. Xu, H., Baracskay, P., O'Neill, J., and Csicsvari, J. (2019). Assembly responses of hippocampal CA1 place cells predict learned behavior in goal-directed spatial tasks on the radial eight-arm maze. *Neuron* 101, 119–132.e4. <https://doi.org/10.1016/j.neuron.2018.11.015>.
60. Kurth-Nelson, Z., Economides, M., Dolan, R.J., and Dayan, P. (2016). Fast sequences of non-spatial state representations in humans. *Neuron* 91, 194–204. <https://doi.org/10.1016/j.neuron.2016.05.028>.
61. Schuck, N.W., Cai, M.B., Wilson, R.C., and Niv, Y. (2016). Human orbitofrontal cortex represents a cognitive map of state space. *Neuron* 91, 1402–1412. <https://doi.org/10.1016/j.neuron.2016.08.019>.
62. Wilson, R.C., Takahashi, Y.K., Schoenbaum, G., and Niv, Y. (2014). Orbitofrontal cortex as a cognitive map of task space. *Neuron* 81, 267–279. <https://doi.org/10.1016/j.neuron.2013.11.005>.
63. Libby, A., and Buschman, T.J. (2021). Rotational dynamics reduce interference between sensory and memory representations. *Nat. Neurosci.* 24, 715–726. <https://doi.org/10.1038/s41593-021-00821-9>.
64. Freedman, D.J., Riesenhuber, M., Poggio, T., and Miller, E.K. (2001). Categorical representation of visual stimuli in the primate prefrontal cortex. *Science* 291, 312–316. <https://doi.org/10.1126/science.291.5502.312>.
65. Freedman, D.J., and Assad, J.A. (2016). Neuronal mechanisms of visual categorization: an abstract view on decision making. *Annu. Rev. Neurosci.* 39, 129–147. <https://doi.org/10.1146/annurev-neuro-071714-033919>.
66. Sarma, A., Masse, N.Y., Wang, X.J., and Freedman, D.J. (2016). Task-specific versus generalized mnemonic representations in parietal and prefrontal cortices. *Nat. Neurosci.* 19, 143–149. <https://doi.org/10.1038/nn.4168>.
67. Fitzgerald, J.K., Freedman, D.J., Fanini, A., Bennur, S., Gold, J.I., and Assad, J.A. (2013). Biased associative representations in parietal cortex. *Neuron* 77, 180–191. <https://doi.org/10.1016/j.neuron.2012.11.014>.
68. Ganguli, S., Bisley, J.W., Roitman, J.D., Shadlen, M.N., Goldberg, M.E., and Miller, K.D. (2008). One-dimensional dynamics of attention and decision making in LIP. *Neuron* 58, 15–25. <https://doi.org/10.1016/j.neuron.2008.01.038>.
69. Duvellé, É., Grieve, R.M., Liu, A., Jedidi-Ayoub, S., Holeniewska, J., Harris, A., Nyberg, N., Donnarumma, F., Lefort, J.M., Jeffery, K.J., et al. (2021). Hippocampal place cells encode global location but not connectivity in a complex space. *Curr. Biol.* 31, 1221–1233.e9.
70. Sutton, R.S., Precup, D., and Singh, S. (1999). Between MDPs and semi-MDPs: a framework for temporal abstraction in reinforcement learning. *Artif. Intell.* 112, 181–211. [https://doi.org/10.1016/S0004-3702\(99\)00052-1](https://doi.org/10.1016/S0004-3702(99)00052-1).
71. Tustison, N.J., Avants, B.B., Cook, P.A., Zheng, Y., Egan, A., Yushkevich, P.A., and Gee, J.C. (2010). N4ITK: improved N3 bias correction. *IEEE Trans. Med. Imaging* 29, 1310–1320. <https://doi.org/10.1109/TMI.2010.2046908>.
72. Klein, A., Ghosh, S.S., Bao, F.S., Giard, J., Häme, Y., Stavsky, E., Lee, N., Rossa, B., Reuter, M., Chaibub Neto, E.C., et al. (2017). Mindboggling morphometry of human brains. *PLoS Comput. Biol.* 13, e1005350. <https://doi.org/10.1371/journal.pcbi.1005350>.
73. Greve, D.N., and Fischl, B. (2009). Accurate and robust brain image alignment using boundary-based registration. *NeuroImage* 48, 63–72. <https://doi.org/10.1016/j.neuroimage.2009.06.060>.
74. Power, J.D., Mitra, A., Laumann, T.O., Snyder, A.Z., Schlaggar, B.L., and Petersen, S.E. (2014). Methods to detect, characterize, and remove motion artifact in resting state fMRI. *NeuroImage* 84, 320–341. <https://doi.org/10.1016/j.neuroimage.2013.08.048>.
75. Fedorenko, E., Duncan, J., and Kanwisher, N. (2013). Broad domain generality in focal regions of frontal and parietal cortex. *Proceedings of the National Academy of Sciences* 110, 16616–16621.
76. Tarhan, L., and Konkle, T. (2020). Reliability-based voxel selection. *NeuroImage* 207, 116350. <https://doi.org/10.1016/j.neuroimage.2019.116350>.



## STAR★METHODS

### KEY RESOURCES TABLE

| REAGENT or RESOURCE   | SOURCE  | IDENTIFIER  |
|---|---|---|
| <b>Deposited data</b>                                       |   |   |
| Human behavioural data                                      | This paper  | <a href="https://osf.io/z25aw/">https://osf.io/z25aw/</a>   |
| Human fMRI data   | This paper  | <a href="https://osf.io/z25aw/">https://osf.io/z25aw/</a>   |
| <b>Software and algorithms</b>                              |   |   |
| Custom code for simulations, experiments, and data analysis | This paper  | doi: <a href="https://zenodo.org/record/8246406">https://zenodo.org/record/8246406</a>  |
| MATLAB 2019   | Mathworks   | <a href="https://uk.mathworks.com/products/matlab.html">https://uk.mathworks.com/products/matlab.html</a> ;<br>RRID: SCR_001622         |
| Python 3.8.5  | Python  | <a href="https://www.python.org/">https://www.python.org/</a> ;<br>RRID: SCR_008394   |
| SPM 12  | Wellcome Centre for Human Neuroimaging, London            | <a href="https://www.fil.ion.ucl.ac.uk/spm/software/spm12/">https://www.fil.ion.ucl.ac.uk/spm/software/spm12/</a> ;<br>RRID: SCR_007037 |
| fMRIPrep 20.2.3   | Poldrack Lab, Stanford                                    | <a href="https://fmriprep.org/en/20.2.3/index.html">https://fmriprep.org/en/20.2.3/index.html</a> ;<br>RRID: SCR_016216                 |
| FSL 5.0.9   | Wellcome Centre for Integrative Neuroimaging, Oxford      | <a href="https://fsl.fmrib.ox.ac.uk/fsl/fslwiki">https://fsl.fmrib.ox.ac.uk/fsl/fslwiki</a> ;<br>RRID: SCR_002823                       |
| AFNI  | NIMH, Bethesda  | <a href="https://afni.nimh.nih.gov/">https://afni.nimh.nih.gov/</a> ;<br>RRID: SCR_005927   |
| FreeSurfer 6.0.1  | Martinos Centre for Biomedical Imaging, Harvard           | <a href="https://surfer.nmr.mgh.harvard.edu/">https://surfer.nmr.mgh.harvard.edu/</a> ;<br>RRID: SCR_001847                             |
| Mindboggle  | Community Project   | <a href="https://mindboggle.info/">https://mindboggle.info/</a> ;<br>RRID: SCR_002438   |
| ANTS 2.2.0  | Brian B. Avants, Nicholas J. Tustison, & Hans J. Johnson, | <a href="https://github.com/ANTsX/ANTs">https://github.com/ANTsX/ANTs</a> ;<br>RRID: SCR_004757   |
| Nipype  | Community Project   | <a href="https://nipype.org/packages/nipype/index.html">https://nipype.org/packages/nipype/index.html</a> ;<br>RRID: SCR_002502         |
| Numpy 1.19.2  | Community Project   | <a href="https://numpy.org/">https://numpy.org/</a> ;<br>RRID: SCR_008633   |
| Scipy 1.6.1   | Community Project   | <a href="https://scipy.org/">https://scipy.org/</a> ;<br>RRID: SCR_008058   |
| Scikit Learn 0.24.2   | Community Project   | <a href="https://scikit-learn.org/stable/">https://scikit-learn.org/stable/</a> ;<br>RRID: SCR_002577                                   |
| Unity Experiment Framework                                  | Jack Brookes  | <a href="https://immersivecognition.com/unity-experiment-framework/">https://immersivecognition.com/unity-experiment-framework/</a>     |

### RESOURCE AVAILABILITY

#### Lead contact

Further information and requests should be directed to and will be fulfilled by the Lead Contact, Paul Muhle-Karbe ([p.muhle-karbe@bham.ac.uk](mailto:p.muhle-karbe@bham.ac.uk)).

#### Materials availability

This study did not generate any new materials.

#### Data and code availability

Code to replicate our analyses is available on Github (doi: <https://zenodo.org/record/8246406>). In line with local ethics guidelines, preprocessed and anonymized group-level human fMRI data is available on OSF (<https://osf.io/z25aw/>).

## METHOD DETAILS

### Participants

Thirty-one human participants were recruited for the experiment through the recruitment system at the Max Planck Institute for Human Development (Berlin). One participant was omitted from the analysis due to a neural structural abnormality and another three participants were omitted due to technical difficulties with the MRI equipment. All analyses were performed on the remaining 27 participants (11 male, 16 female; age  $27.3 \pm 4.4$  years). Participants were compensated for their time at a base rate of €10/hour, plus an extra €10 for participating in an MRI experiment, and finally an additional bonus of up to 10€ (5€ per session) depending on their performance. Informed consent was given before the start of the experiment. The study was approved by the Department of Education and Psychology at the Freie Universität Berlin and the Medical Science Inter-Divisional Research Ethics Committee (R49432/RE001) at the University of Oxford.

### Design, task and procedures

The experiment involved two sessions undertaken on different days. In both sessions, participants performed a computerised task that was built and delivered in the Unity 3D games environment. The task involved navigating an avatar through a grid world to collect rewards. On day 1, participants performed a training task outside of the scanner, using the arrow keys on a laptop computer to move the avatar through the environment (see Figure 1A). On day 2 ( $32.0 \pm 3.6$  hours later), they performed the task lying supine in an MRI scanner, viewing the screen through a mirror and using an MRI-compatible button box to respond.

On both days, the grid world environment was composed of four adjoining rooms arranged in a square.<sup>69</sup> We refer to the rooms as southwest (SW), northwest (NW), southeast (SE) and northeast (NE) rooms. Each room was composed of  $4 \times 4$  grid squares and was connected to the two cardinal adjacent rooms (e.g., SW was connected to SE and NW but not NE) via a single “bridge” square. It thus mirrored the classic “four rooms” environment commonly used in AI research.<sup>70</sup> At each point in the trial, participants could only see the  $4 \times 4$  squares of the currently occupied room, plus the two additional bridge squares; the other rooms were offscreen. One square of each room contained a boulder, and two of the four boulders in the environment were associated with a reward (the reward was revealed when the avatar collided with the boulder). During training, the grid squares were differently coloured in each of the four rooms (in order to help people learn to navigate); during test, they were all purple. Traversing a bridge square incurred a variable delay during which the full map was briefly shown. This was to encourage participants to consider their room choices carefully before moving between rooms, and later on day 2, to more easily separate the BOLD response pertaining to the occupation of different rooms.

On both days, the task was divided into blocks of 16 trials ( $n = 4$  during day 1;  $n = 6$  during day 2). In the scanner, these constituted 6 independent scanner runs. On each trial, participants began in the inner corner of a randomly chosen room (they could identify the room by the locations of the visible bridge squares). Before navigation began, participants were shown a contextual cue, which was a picture of one of two food items (Figure 1C). Unbeknownst to participants, each cue disclosed one reward location conditional on the other for that trial. For example, in scanner run 1, cue A (a martini icon) indicated that the rewards were in rooms lying in the same horizontal axis, and cue B (a peanut icon) that the rewards were in the same vertical axis (with neither disclosing which specific rooms). Different pairs of food items were chosen on day 1, and then on blocks 1-2, 3-4, and 5-6 of day 2 (4 pairs total) so that participants had to generalise the reward co-variance structure to previously unseen items. Food icons used for the day 2 (scanner) task included watermelon, cheese, mushroom, avocado, pineapple and banana.

Participants navigated freely (up, down, left, right) using buttons, causing the avatar to move within the environment (the background grid remaining fixed). When participants alighted on a boulder that was associated with a reward, the reward was revealed by showing the food item that had been cued on that trial, before navigation could recommence. Participants were instructed to find the two rewards as quickly as possible and received a *trial score* that was equal to the number of seconds the participant had remaining on their timer at the end of each trial. If participants took longer than the timer deadline to find both rewards, 20 points were deducted from their total trial score. The task was calibrated so that a participant that ignored the cues and navigated to boulders in any order would only meet the deadline on approximately 50% of trials, set to be 40s on day 1 and 50s in the scanner on day 2. Aggregate trial score was converted to a financial bonus at the end of the experiment.

The timing of events within each trial were as follows. Each trial started with the controls disabled, and the location of the avatar in the start room was shown for 2.5s. The contextual cue was then displayed enlarged in the centre of the screen for 1.5s. After a further 1s the controls were enabled, and participants were able to move the avatar through the environment by pressing arrow keys (day 1) or button box keys (day 2). At this point, the timer (visible in the top right hand corner of the screen) started ticking down from a deadline value (40s on day 1, 50s on day 2). On day 2, participants could move the avatar at a maximum speed of 1 grid square every 0.4s (increased from 0.25s on day 1, to ensure participants remained in each room long enough to obtain a clear per-room neural signal). When moving through a hallway, controls were disabled for a period of time before players were able to move again, where this period was drawn from a truncated exponential distribution (mean 2s; min 1.5s; max 7s). When the avatar collided with a boulder, the controls were again disabled for a period (sampled from truncated exponential with mean 2s; min 1s; max 5s) while either a reward or no reward was shown. At the end of the trial, a message saying “well done” appeared on the screen and the participant’s total score was visibly updated using the remaining seconds left on the timer, which corresponded to additional points. After each trial, participants were shown a black screen for a period before the next trial began (ITI sampled from truncated exponential distribution with mean 2.5s, min 1.5s, max 7s).

In our design, we consider the cues to be “contexts” signaling whether rewards were found on the horizontal or vertical axes of the four rooms environment. However, because rewards could be in any room (i.e., in a vertical condition they could be in the SW and NW or the SE and NE), occupancy probabilities were closely matched across contexts. Nevertheless, to ensure good coverage of the environment, and to attempt to match transitions as well as occupancy, we also introduced a “robot control” phase in every trial, in which either the robot or the participant began by controlling the avatar. If the robot controlled the avatar, it moved at approximately the same pace as an average participant, and typically to a non-rewarded room, where it made a beeline for the boulder. Every time the computer (participant) reached a boulder, the control was passed back to the participant (computer) until the next boulder was reached. On average, the amount of time spent per trial under control of the robot was 11.7 seconds, compared to 14.3 seconds under participant control. We also verified that the neural representation of space independently for the human- and robot-controlled phase (Figure S6D-S6E).

The behavioural training session (day 1) began with two practice trials, which used different contextual cues, which did not signal the location of one cue conditional on the other. During the training session (day 1), cues were blocked, so that participants alternated between horizontal and vertical contexts in an ABAB design. During day 2, in the scanner, contexts were interleaved from trial to trial, so that different contexts were not associated with distinct, prolonged temporal episodes.

Each run contained 16 trials, which were balanced across pairs of runs with the same reward cues. Trials were balanced across the two cues (32 trials / 2 cues = 16), starting rooms (16 trials per cue / 4 rooms = 4), whether the start room was rewarded or not (4 trials per room per cue / 2 = 2), and whether the participant foraged first or the robot foraged first (2 trials rewarded per room per cue / 2 = 1). Trial ordering was randomised across participants.

At the end of the scanning sessions, participants responded to a situational quiz which examined their explicit understanding of the reward covariance rules. They were asked four questions of the form “*You have just found a cheese in the top right room. Which room will the other cheese be in?*”, and four questions which tested their counterfactual understanding, such as “*You were looking for a cheese and did NOT find one in the bottom right room. Which rooms will contain the two cheeses?*”. Each participant was assessed with a version of the quiz that mentioned the reward pair that they had most recently observed (those in the final two runs). The maximum possible score was 8. We examined the correlation across participants between these quiz scores and (1) the average first room choice accuracy in the scanning session, and (2) the mean trial score in the scanning session.

### Behavioural analysis

We computed three behavioural metrics. Firstly, we considered the trial score, which is roughly proportional to the average time taken to complete a trial. Secondly, we computed the *transition bias*, which is the relative fraction of transitions between horizontal and vertical rooms in the appropriate context (H or V), computed across both human- and computer-controlled events:

$$p(\text{horizontal}|\text{H}) + p(\text{vertical}|\text{V}) - p(\text{horizontal}|\text{V}) - p(\text{vertical}|\text{H})$$

Thirdly, we computed first choice accuracy, which is the probability that participants made an optimal transition from the first room occupied to the second. This measure is particularly sensitive because a participant who perfectly understand the meaning of the contextual cues can always head to the correct 2<sup>nd</sup> room, on the basis of whether the 1<sup>st</sup> (starting) room contains a reward or not. First choice accuracy was quite highly correlated with transition bias ( $r = 0.86$ ).

### Computational model

We defined a simulated environment corresponding to the four rooms arena, in which locations were denoted by values in the range  $[-1, 1]$  in both the  $x$  and the  $y$  dimension. We rescaled participants’ observed movement trajectories through the grid environment so that they mapped onto this simulated environment and located the boulders at their approximately corresponding positions. This allowed us to model neural responses using an encoding model that consisted of simulated place cells. The place cells exhibited bivariate Gaussian response fields that regularly tiled the arena on a  $10 \times 10$  square lattice (but the results we describe were very similar, albeit but more variable, if we drew their tuning preferences from random uniform distributions; we also verified that almost identical results are obtained if we truncate place fields so that they do not straddle different rooms). We assume that there are 200 place cells in each context (two cells coding for each location).

In this simple model, we define the place field of neuron  $i$  in context  $c$  as peaking at an  $[x, y]$  location  $\theta_i(c)$ . We note that the vector of place fields in the two contexts  $\theta(c = V)$  and  $\theta(c = H)$  may be the same, or partially or fully different. This is controlled by the parameter  $\beta$ , which determines the fraction of cells for which  $\theta(c = H) \neq \theta(c = V)$ . Thus if  $\beta = 0$  then all cells code for the same location regardless of context (no remapping), if  $\beta = 0.5$  then 50/100 cells exhibit overlapping place fields between contexts (partial remapping), and if  $\beta = 1$  then all cells change their tuning between contexts (full remapping).

Thus, in any given context we can estimate the neural response of neuron  $i$  on time step  $t$  as being

$$R_{t,i} = |\omega| \times f(s_t|\theta_i(c), \Sigma) + (1 - |\omega|) \times f(g_t|\theta_i(c), \Sigma) + \gamma \times h(c)$$

In the expression above,  $f(\cdot|\theta_i, \sigma)$  is the bivariate normal distribution evaluated at the preferred  $[x, y]$  tuning location (place field) for neuron  $i$  in context  $c$ , and  $\Sigma = [0.25, 0; 0 \ 0.25]$ . We define the current  $[x, y]$  location of the agent as  $s_t$ , whereas  $g_t$  is the  $[x, y]$  coordinate location of the goal to which the agent is headed on the current timestep. Finally,  $h(c)$  is a context-specific neural signal, which depends uniquely on the current context (H or V) active on that timestep and not on the location of the agent or goal.

In addition to orthogonalization (controlled by  $\beta$ ), separation and compression are controlled via two further free parameters. The gain parameter  $\gamma$  determines the relative influence of  $h(c)$ , the (place-insensitive) signal coding for context, which has the effect of neural separation between neural manifolds for space in each context. Finally, the mixing parameter  $\omega$  determines the relative influence of the current ( $s_t$ ) and prospective ( $g_t$ ) location on the neural population response. Where  $\omega = 0$ , only the participant's current location is encoded, as in a "classical" place field model. Where  $\omega > 0$ , the model encodes a mixture of the current location and the prospective (goal) location. We also allow for  $\omega < 0$ ; in this case, we use the closely related expression

$$R_{tj} = |\omega| \times f(s_t|\theta_i(c), \Sigma) + (1 - |\omega|) \times f(g'_t|\theta_i(c), \Sigma) + \gamma \times h(c)$$

Where  $g'_t$  is a fictitious goal location that is swapped on the horizontal / vertical axis, as if participants were prospectively encoding horizontal locations in the V conditions and vertical locations in the H condition. This entails that space is compressed along the dimension perpendicular to the axis on which the two goals can be found. We call this anti-compression.

We use the observed individual trajectories, and the actual prospective goals (i.e., where participants were genuinely headed) at each time point to evaluate the model for each agent. This provides us with a neural response matrix of size  $200 \times t$  in each of 48 trials in context H and 48 trials in context V. We then average those timepoints in which the avatar was in the SW, NW, SE and NE rooms for each context, yielding a  $200 \times 8$  matrix, which we use to generate an  $8 \times 8$  RDM (expressing correlation distance) for each simulated participant. For visualisation, we average these RDMs across the simulated cohort, and plot them using multidimensional scaling (e.g., [Figure 2C](#)).

## fMRI data collection and pre-processing

### Anatomical MRI data

MRI data were acquired at the Max Planck Institute for Human Development in Berlin using a 32-channel head coil on a 3T Siemens Magnetom TrioTrio MRI scanner (Siemens, Erlangen, Germany). At the start of the scanning session, a T1-weighted (T1w) high-resolution anatomical image was obtained using a Magnetization Prepared Rapid Gradient Echo (MPRAGE) sequences (sequence parameters: repetition time (TR) = 2500 ms, echo time (TE) = 4.77 ms, flip angle =  $7^\circ$ , field of view (FOV) = 256 mm; voxel size =  $1 \times 1 \times 1$  mm).

Data were processed within the fMRIPrep framework. The T1w image was corrected for intensity non-uniformity with N4BiasField-Correction,<sup>71</sup> distributed with ANTs 2.2.0, and used as T1w-reference throughout the workflow. The T1w-reference was then skull-stripped with a Nipype implementation of the antsBrainExtraction.sh workflow (from ANTs), using OASIS30ANTs as target template. Brain tissue segmentation of cerebrospinal fluid, white matter and gray matter was performed on the brain-extracted T1w using fast (FSL 5.0.9). Brain surfaces were reconstructed using recon-all (FreeSurfer 6.0.1), and the brain mask estimated previously was refined with a custom variation of the method to reconcile ANTs-derived and FreeSurfer-derived segmentations of the cortical gray matter of Mindboggle.<sup>72</sup> Volume-based spatial normalization to MNI space (MNI152NLin2009cAsym) was performed through nonlinear registration with antsRegistration (ANTs 2.2.0), using brain-extracted versions of both T1w reference and the T1w template.

### Functional MRI data

Functional MRI data were acquired using a T2-weighted (T2w) echo planar imaging (EPI) pulse sequence sensitive to BOLD contrast (sequences parameters: TR = 2000 ms, TE = 30 ms, FOV =  $192 \times 192$  mm, flip angle =  $80^\circ$ , voxel size =  $3 \times 3 \times 3$  mm). The task was divided into 6 functional runs, each lasting between 10 and 15 minutes, depending on participant performance.

For each of the six scanning runs, the following pre-processing steps were performed. Initially, a reference volume and its skull-stripped version were generated using a custom methodology of fMRIPrep. A B0-nonuniformity map was estimated based on two EPI references with opposing phase-encoding directions, with 3dQwarp (AFNI 20160207). Based on the estimated susceptibility distortion, a corrected EPI reference was calculated for a more accurate co-registration with the anatomical reference. The BOLD reference was then co-registered to the T1w reference using bregister (FreeSurfer) which implements boundary-based registration.<sup>73</sup> Co-registration was configured with six degrees of freedom. Head-motion parameters with respect to the BOLD reference (transformation matrices, and six corresponding rotation and translation parameters) are estimated before any spatiotemporal filtering using mcflirt (FSL 5.0.9). BOLD runs were slice-time corrected using 3dTshift from AFNI 20160207 and the time-series were resampled to their original native space as well as to the standard MNI space. BOLD data were moreover smoothed with a 6 mm full-width half-maximum Gaussian kernel. A reference volume and its skull-stripped version were generated using a custom methodology of fMRIPrep. Several confounding time-series were calculated based on the pre-processed BOLD: framewise displacement, DVARS and three region-wise global signals. FD and DVARS are calculated for each functional run, both using their implementations in Nipype (following the definitions by Power et al.<sup>74</sup>).

### fMRI analysis: first level GLMs and ROI definition

The pre-processed BOLD timeseries data were modelled with general linear models (GLMs) that contained regressors for different task events. The first GLM (GLM1) contained the following regressors: the contextual cue, the movement period of each trial before the first feedback when subjects had no knowledge about the locations of the reward, the subsequent movement periods when the agent was occupying a room without a reward (pre-goal room period), the subsequent movement periods when the agent was occupying a room with a reward (goal room period), those periods when the agent occupied a hallway between adjacent rooms, the feedback periods when a reward was presented, and the feedback periods when no reward was presented. Note that all movement



periods after the first feedback were labelled in a consistent manner, based on the presence vs. absence of reward in the current room. Hence, the second movement period of each trial (i.e., when the agent is still in the first room after receiving the first feedback) would be treated as pre-goal room or goal room, depending on the outcome of the first feedback.

We modelled data from both the self-directed and robot control periods together to ensure adequate coverage of space in both contexts (see [Figures S6D](#) and [S6E](#) for a control analysis showing aligned spatial representations during both types of events). Note that we defined separate regressors for pre-goal room periods and goal room periods for each of the four rooms of the grid world (SW, NW, NE, SE) and each behavioural context (vertical vs. horizontal goal alignment), resulting in 16 regressors for movement periods after the first feedback. The GLM also contained nuisance regressors pertaining to participants' head motion (three rotation parameters and three translation parameters), the global signal in the white matter, and the framewise displacement. All regressors were modelled with variable durations from start to finish and convolved with a canonical haemodynamic response function. To ensure sufficient trial counts, we concatenated the BOLD time-series data across odd and even scanning runs (using the `spm_fmri_concatenate` function). We conducted a second GLM (GLM2) that defined movement periods in terms of the current location of the agent (SW, NW, NE, SE) and the prospective location of the navigational goal (SW, NW, NE, SE). This GLM collapsed events across behavioural contexts (i.e., trials with vertical vs. horizontal goal alignment), thereby also resulting in 16 regressors for movement periods after the first feedback.

Five regions-of-interest (ROIs) were defined based on existing atlases. We used the Wake Forest University Pickatlas (integrated into SPM) to define ROIs for the hippocampus (bilateral areas labelled *Hippocampus*), orbitofrontal cortex (bilateral areas labelled *Frontal\_Inf\_Orb*; *Frontal\_Mid\_Orb*; *Frontal\_Sup\_Orb*), and visual cortex (bilateral areas labelled *Occipital\_Mid*). ROIs for prefrontal and posterior parietal cortices were defined based on an atlas provided by Fedorenko et al.<sup>75</sup> that delineates frontoparietal brain areas implicated in cognitive control across a variety of cognitive domains (the whole atlas is available for download at <http://imaging.mrc-cbu.cam.ac.uk/imaging/MDsystem>).

### fMRI analysis: neural geometry

To compute the neural geometry, we obtained the multivariate pattern evoked by each of the 16 predictors (in either GLM1 or GLM2) for each participant in each scanner run. Thus, for each region (or searchlight) this yielded a  $n(v_r) \times 16 \times 6$  data array, where  $n(v_r)$  is the number of voxels in region  $r$ . We collapsed over odd and even runs, giving us two  $n(v_r) \times 16$  arrays. Using a previously described method called 'reliability-based voxel selection'<sup>76</sup>, we began by identifying eligible voxels for multivariate analysis (feature selection). We correlated, in each individual voxel, the pattern of activity over the 16 conditions between odd and even runs. Voxels with Spearman's  $r > 0$  were included in all multivariate analysis. This left a minimum of 1470, 1931, 709, 320 and 1335 voxels in visual, PPC, PFC, hippocampus and OFC ROIs respectively. This feature selection method ensures that only those voxels with consistent patterns across runs (i.e., those with higher signal) are included in the analysis (but it does not specify what the pattern should be in these voxels). To verify that this did not bias our analysis in any way, we reran all analyses in the paper on shuffled data to which we applied the same feature selection methods. To achieve this, we shuffled the mapping across voxels independently between training and test, creating a dataset with equivalent summary statistics but no train-test consistency, and reran our analyses including the feature selection stage. We observed no deviation from the expected null distribution in this case. Next, we used singular value decomposition to reduce the dimensionality to  $d$  dimension within the ROI; for all statistical analyses, we used  $d = 10$  (we did not apply this step for MDS visualisation). The first 10 principal components captured about ~60% of the variance in most regions and participants.

At this stage (where required, i.e., for [Figures 3](#) and [4](#)) we separated the ROIs into pre-goal room and goal room periods. This was desirable because of the large offset in behaviour between these periods, but we obtain very similar results when we analyse all the data together ([Figure 5](#)). In each case, we computed RDMs ( $8 \times 8$  or  $16 \times 16$ ) in cross-validation; this means that we computed  $RDM_{ij}$  which is the dissimilarity between the neural pattern from the  $i^{\text{th}}$  condition in odd scanner runs to the  $j^{\text{th}}$  condition in even scanner runs. We then averaged these RDMs about the diagonal and regressed the lower triangle of the RDM against that of a predictor matrix composed of one or more (standardised) model RDMs, obtaining beta coefficients for their competitive fit. The diagonal is not zero for our data RDMs because of the cross-validation step, so we set it to zero for MDS visualisation only; note that this has no impact on our statistical analysis. The RDMs we generated were designed to be as orthogonal as possible. For example, the "map" RDM captures similarity structure over and above that in the "room" RDM by predicting larger distances on the diagonals (e.g. NE to SW) than edges (e.g. NE to NW; see [Figure 3E](#)). The median Pearson correlation between predictors was  $r = 0.175$  and no pair of predictors had a Pearson's correlation that exceeded 0.5.

For the "scores" analysis, we generated three perfectly orthogonal matrices that we call *compression*, *separation* and *map*. Each of these score matrices is the same size as the RDM and comprises binary values (+1 and -1) for key condition pairs; each sums to zero. Each score matrix is multiplied elementwise with the data RDM for each participant and averaged, yielding a score that is  $> 0$  if matrix values set to 1 are more dissimilar than those set to -1 and  $< 0$  otherwise. This allows us to do group-level one-sample t-tests (against zero) to test for these three effects.

The *compression* matrix tests whether the east and west rooms are more similar in the H condition, and the north and south rooms in the V condition, than the converse (it is thus the subtraction of the compression and anti-compression matrices shown in [Figure 3E](#)). The *separation* matrix sets values between contexts to +1 and those within contexts to -1, excluding the minor diagonals (i.e., the dissimilarity between each room and itself across contexts). The *map* matrix tests whether each plane is shaped like a quadrilateral, mirroring the geometry of the four rooms environment. To this end, the scores matrix has values of +1 for the diagonals (e.g., SE to

NW) and -1 for the cardinals (e.g., SE to NE; this is done in subsets to ensure the matrix sums to zero). The resulting vector of scores across the participant cohort is correlated with behavioural measures, including first choice accuracy and transition bias, using Pearson's correlation.

In the angle analysis, we assess the angle between the north-south and east-west vectors within contexts, across contexts, and across contexts and goal room period. We do this using the full  $16 \times 16$  matrix. We first compute the difference in high-dimensional vector coding for each room in each context and period, leading to a data matrix of size  $n(v_r) \times 16 \times 16$ . We manually compute the angle between those edges predicted to be parallel, orthogonal or inverted in the model in [Figure 5C](#) and plot these in [Figure 5D](#).

To detect potential signals outside of the chosen ROIs, we repeated the "scores" analysis, described above, in a whole-brain searchlight approach, where RSA was conducted at each voxel with a group of surrounding voxels (spherical searchlight radius = 12 mm). Analogously to the ROI analyses, we conducted separate analyses for pre-goal room period and goal room period. For each voxel within the searchlight, we extracted the 8 beta coefficients from the GLM corresponding to the regressors for each room (SW, NW, NE, SE) and context (vertical, horizontal). We next applied voxel selection and dimensionality reduction, as described above, and computed cross-validated RDMs ( $8 \times 8$ ), which were multiplied elementwise with each predictor matrix, yielding three whole-brain maps with regression coefficients for each subject. These maps were smoothed using an 8 mm FWHM Gaussian kernel. Statistical significance was established separately for each voxel by testing the regression coefficients against zero using one-sample t-tests. Correction for multiple comparison was conducted via family-wise error correction ( $p < 0.05$ ) as implemented in SPM 12.