



Published in final edited form as:

Nat Neurosci. 2017 April ; 20(4): 590–601. doi:10.1038/nn.4509.

Persistently active neurons in human medial frontal and medial temporal lobe support working memory

J Kami ski^{1,4}, S Sullivan¹, JM Chung², IB Ross³, AN Mamelak¹, and U Rutishauser^{1,2,4}

¹Department of Neurosurgery, Cedars-Sinai Medical Center, Los Angeles, California, USA

²Department of Neurology, Cedars-Sinai Medical Center, Los Angeles, California, USA

³Department of Neurosurgery, Huntington Memorial Hospital, Pasadena, California, USA

⁴Division of Biology and Biological Engineering, California Institute of Technology, Pasadena, CA

Abstract

Persistent neural activity is a putative mechanism for the maintenance of working memories. Persistent activity relies on the activity of a distributed network of areas, but the differential contribution of each area remains unclear. We recorded single neurons in the human medial frontal cortex and the medial temporal lobe while subjects held up to three items in memory. We found persistently active neurons in both areas. Persistent activity of hippocampal and amygdala neurons was stimulus-specific, formed stable attractors, and was predictive of memory content. Medial frontal cortex persistent activity, on the other hand, was modulated by memory load and task set but was not stimulus-specific. Trial-by-trial variability in persistent activity in both areas was related to memory strength, because it predicted the speed and accuracy by which stimuli were remembered. This work reveals, in humans, direct evidence for a distributed network of persistently active neurons supporting working memory maintenance.

Introduction

The ability to store information in an active and readily available state for short periods of time is fundamental for cognition. Working memory (WM) is essential for many high level cognitive skills, including inference, decision making, mental calculations, and awareness¹. The prevailing models of WM posit that in the absence of external stimuli, memoranda are maintained by persistent neuronal activity^{2,3}. Signatures of sustained activity have been observed in a variety of brain areas in macaques including dorsolateral prefrontal cortex (dlPFC)^{4–6}, parietal cortex⁷, inferior temporal (IT) cortex^{8,9}, Entorhinal Cortex¹⁰, and

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:http://www.nature.com/authors/editorial_policies/license.html#terms

Corresponding author: ueli.rutishauser@cshs.org.

Contributions

J.K. and U.R. designed the experiments. J.K. and U.R. performed experiments. J.K., S.S. and U.R. performed analysis. A.N.M. and I.B.R. performed surgery. J.M.C. provided patient care. J.K. and U.R. wrote the paper. All of the authors discussed the results at all stages of the project.

Competing financial interests

The authors declare no competing financial interests.

medial frontal cortex (MFC) ¹¹. It is thus thought that WM relies on different types of persistent activity provided by a distributed set of brain areas. Indeed, some brain areas in non-human primates exhibit stimulus-selective ^{4,5,7-9} persistent activity, whereas others do not ¹². However, the relationship and importance of these different kinds of persistent activity for human WM remain unclear.

In humans, a number of brain areas are thought to be essential for forming, maintaining and retrieving short-term memories ¹³. Here, our focus is on two of these areas: the MTL and the MFC. In neuroimaging studies MFC activity is consistently related to WM ¹⁴ and it has been suggested that the MFC supports executive functions through persistent neuronal activity ¹⁵. However, there is at present no direct neuronal evidence of such activity in the human MFC. The MTL, on the other hand, is traditionally thought to be required for the formation of new long-term memories but not for WM ¹⁶. However, it has become clear that the MTL also plays an important role for WM ¹⁷ whenever subjects are distracted or are maintaining an amount of information close to their WM capacity. Indeed, studies of patients with damage to the MTL reveal deficits in retaining information for more than a few seconds in WM tasks under some circumstances ¹⁷. Also, studies utilizing intracranial recordings have begun to reveal field potential signatures of persistent activity during maintenance of WM in the MTL ¹⁸ and recordings from human lateral temporal cortex have revealed elevated activity of neurons during encoding and retrieval of short term memories ¹⁹. However, there is presently no direct single-neuron evidence of persistently active MTL or MFC neurons during WM maintenance in humans.

Attractors are dynamically stable patterns of neuronal activity that have been an influential framework for conceptualizing how memories are maintained by persistent activity ^{20,21}. In this framework, the brain has large numbers of different attractors, each of which corresponds to a different specific memory. During maintenance, at most, a small subset of these attractors can be active simultaneously. Consequently, the extent to which the pattern of persistent activity during a given trial remains close to the attractor(s) should predict the quality of the memory trace as assessed by later behavior. Indeed, in macaques, the extent of drift of neural activity away from attractors formed by persisted activity is predictive of error magnitude in recalled spatial location ²². While a major concept in theoretical work ^{20,21,23}, there is little direct evidence for the existence and importance of attractors for WM maintenance in humans.

Here, we used human single-neuron recordings performed in neurosurgical patients to investigate the properties of persistent neuronal activity and its relationship to behavior. We recorded simultaneously from two areas each in the Medial Temporal Lobe (MTL, hippocampus and amygdala) and the Medial Frontal Cortex (MFC; we recorded from Dorsal Anterior Cingulate Cortex, dACC and pre-Supplementary Motor Area, pre-SMA). As a starting point, we screened for highly selective and sparsely responsive “concept cells”. Such cells are characterized by a highly selective and invariant response for abstract concepts, such as a particular person or building ^{24,25}. Here, we test whether concept cells exhibit persistent activity and whether such activity is related to WM. As in previous work ²⁴, we used a screening task to select the images with the best response before our actual

experiment. We then used only the images for which selective neurons were present in the subsequent WM task.

Results

Task and behavior

Subjects (14 sessions from 13 patients) performed a working memory (WM) task (“Sternberg task”). In each trial, subjects were asked to memorize (“encode”) 1–3 images that were presented sequentially for 1s each. After a waiting (maintenance) period of at least 2.5 s and at most 2.8s, subjects were asked to judge whether a probe stimulus was identical to one of the 1–3 images held in memory (Fig. 1a). We refer to the number of images held in memory as “load” throughout. Subjects pressed either the “Yes” or “No” button to provide their answers (button identity was reversed in the middle of the experiment). Patients performed well: average accuracy was $89.55\% \pm 5.66\%$ (Fig. 1c), and median reaction time (for correct trials only) was 1.09 ± 0.28 s (\pm sd). For all subsequent analysis, we used only correct trials unless indicated otherwise. Reaction time (RT) increased as a function of load (Fig. 1d, $F[2,26]=4.11$; $P=0.025$, permuted repeated-measures ANOVA). Also, RT was significantly faster in response to probes that were held in memory (IN= 1.03 ± 0.29 sec) compared to probes not held in memory (OUT= 1.2 ± 0.3 sec; permuted paired t-test: $t[13]=3.35$; $p=0.005$). Moreover, correct in/out decisions were made more quickly than incorrect decisions (1.09 ± 0.28 s vs. 1.52 ± 0.74 s, permuted paired t-test: $t[13]=3.81$; $p=0.0005$). Together, this shows that our patients performed the task accurately and exhibited the expected patterns of RT differences²⁶.

Screening task

We customized the set of images to be memorized for each patient. Each patient performed a screening task 2–3 hours prior to the experiment. During this task, 54–64 images were presented in random order 6 times for 1s each. We then processed the data to identify the 5 images with the best image-selective responses (see methods; see Supplementary Fig. 2 for examples) which were then used for the WM experiment. We will refer to these images as Image A-E throughout the manuscript (A,B,C,D, or E was set arbitrarily but fixed for each patient).

Electrophysiology

Across all patients, we isolated 651 (47 ± 22.1 per session) putative single units from the dorsal Anterior Cingulate Cortex (dACC; $N=120$), pre-Supplementary Motor Area (pre-SMA; $N=184$), amygdala (Amg; $N=195$) and hippocampus (Hipp; $N=152$, Fig. 1b shows the locations of recording sites in MNI space). Spike sorting quality was assessed quantitatively (Supplementary Fig. 1). Throughout the manuscript, we use the terms neuron and cell interchangeably to refer to a putative single unit.

Classes of neuronal responses

We selected for and quantified the properties of three groups of neurons: i) visually selective concept cells, ii) maintenance neurons that exhibit elevated activity in the absence of visual stimuli, and iii) probe neurons, which responded only during retrieval. Together, these three

groups of cells provide a comprehensive inventory of cellular activity during WM encoding, maintenance, and retrieval.

Identification of concept cells

We identified concept cells²⁴ by testing whether their firing rate significantly co-varied as a function of picture identity (using a permuted one-way ANOVA with x groups, where x is the number of unique images followed by a permuted t-test for the image with the maximal response against all other images). In the initial screening sessions, 99 cells (out of 670, 14.77%) qualified as concept cells. In the Sternberg task, 93 cells (out of 465, 14.2%) qualified as concept cells (Fig. 2a shows an example). Most concept cells were located either in the amygdala (Fig. 2b, 25.4% of all amygdala cells, Permutation test, $P=0.002$ vs. chance) or the hippocampus (12.1%, Permutation test, $P=0.002$ vs. chance). The number of concept cells was not larger than expected by chance in the dACC (4.35%, Permutation test, $P=0.21$) and pre-SMA (2.55%, Permutation test, $P=0.806$). We will use the term preferred image for the picture for which a given cell is selective and the term non-preferred images for all other pictures. The response of concept cells was highly selective: the average depth of selectivity (S) index was 0.68 (see Supplementary Fig. 2b and methods). During the screening task, only $9.61\% \pm 9.2\%$ (\pm s.d.) of the presented pictures evoked a selective response in at least one cell. In contrast, during the Sternberg task, on average $56.92\% \pm 36.37\%$ (\pm s.d., permuted paired t-test, $t[12]=5.48$, $P=0.001$) of all presented pictures elicited a response in at least one concept cell (Fig. 2c). 86% of the units which were identified as significantly selective in the screening task remained selective for the same visual stimulus during the Sternberg task. This indicates that the screening procedure was effective, that recordings were stable, and that most neurons maintained their tuning between sessions.

Activity of concept cells during working memory maintenance

We found that concept cells located in the MTL (hippocampus and amygdala) remained active during memory maintenance if the preferred image of a cell was presently held in memory (Fig. 3a–d, Fig. 4a, Supplementary Fig. 3a–b). Note that the preferred and non-preferred images for each cell were defined based only on the activity observed during encoding, making this analysis unbiased. We used the picture selectivity index (PSI) to quantify the activity of each concept cell during maintenance (see methods). During maintenance, the PSI for concept cells recorded in the amygdala ($n=57$) was 49.2 ± 65.88 (permuted paired t-test, $t(56)=4.79$, $P=0.0005$ vs. zero, Fig. 4b). Similarly, the PSI for concept cells recorded in the hippocampus ($n=21$) was 28.8 ± 51.85 (permuted paired t-test, $t(20)=2.54$, $P=0.0025$ vs. zero). Tested individually for each neuron, 44.9% ($N=35$) of the concept cells in the MTL exhibited persistent activity during WM maintenance ($p<0.05$, one-tailed test). This suggests that persistent activity of concept cells in the MTL supports WM. If so, we expected that the PSI systematically co-varied with memory load and success/failure in maintaining a memory. We next investigated this hypothesis.

We found that the higher the load, the lower the PSI (permuted repeated-measures 1×3 ANOVA, $F(2,112)=7.31$, $P=0.0005$ for amygdala and $F(2,40)=8.48$, $P=0.0005$ for hippocampus, Fig. 4c). Amygdala neurons had significant PSI values for all loads, whereas hippocampal neurons maintained above chance values only for loads 1–2. We next tested

whether the drop in PSI as a function of load could be explained by lingering visually-evoked activity. If so, persistent activity should only be present for the picture encoded last⁹. However, we found no significant relationship between the position during encoding (using loads 2–3 only) and the PSI in load 2 (permuted paired t-test, $t(77)=0.38$, $P=0.702$) nor in load 3 (permuted repeated-measures 1×3 ANOVA, $F[2,231]=0.135$, $P=0.873$; Fig. 4d, Fig.3c and Fig.3d). Thus, the persistent activity during maintenance could not be explained by lingering visually-evoked activity.

We next tested whether the activity of concept cells during the maintenance period correlated with the correctness of the later response. For the subset of trials during which the preferred stimulus of a concept cell was held in memory, persistent activity of concept cells was significantly larger in correct compared to incorrect trials (permuted paired t-test, $t(77)=4.92$, $P=0.005$; see Fig. 4e, left side). In contrast, when the preferred stimulus of a cell was not held in memory, there was no significant difference between correct and incorrect trials (Fig. 4e, permuted paired t-test, $t(77)=0.76$, $P=0.48$).

Lastly, if cells show persistent activity, it should be possible to select for concept cells using only maintenance-activity (load 1 only). Indeed, we found that in both the Amygdala (14.9%, $N=29$, Permutation test, $P=0.002$) and the Hippocampus (11.8% $N=18$, Permutation test, $P=0.002$), a significant proportion of neurons qualified as concept cells using this approach.

Together, this shows that concept cells exhibit selective persistent activity during maintenance of items in WM and that the amplitude of this activity was correlated with the quality of the memory trace.

Maintenance neurons

The second group of neurons we characterized were maintenance neurons. These neurons increased their firing rate relative to baseline during the maintenance period regardless of the stimulus that was held in memory. While we observed maintenance neurons in all recorded areas (Fig. 5c), they were most prominent in the two areas of the MFC ($\chi^2[1]=25.08$; $P=5.473e-6$): pre-SMA and dACC, in which 31.3% and 21.2% of all recorded neurons were maintenance neurons, respectively. Was the activity of maintenance neurons predictive of which image(s) were held in memory? Using the same approach as in the case of concept cells (see methods) we found that the activity of maintenance neurons was not indicative of the images held in memory; only 5 (3.97%) neurons out of 126 showed significant differences ($P=0.829$ against scrambled data).

We next tested whether the firing rate of maintenance neurons was associated with load and retrieval performance (RT and accuracy). Indeed, during maintenance, the firing rates of subsets of maintenance neurons in dACC (10.7%) and pre-SMA (46.3%) co-varied systematically with load (Fig. 5d, Supplementary Fig. 3c). Also, in pre-SMA, the majority of such neurons decreased their firing rate as a function of load (9.2% vs 37%, Fig. 5d and Fig. 5a). In addition, the activity of 18.5% of maintenance neurons recorded in the pre-SMA and 10.3% (3 cells) of amygdala maintenance neurons differentiated between slow and fast later memory retrieval (median split computed independently for each load and in and out

conditions, median was computed individually for every subject to account for individual differences in RT; Fig. 5e, Supplementary Fig. 3d). Lastly, the mean firing rate of maintenance neurons in both dACC and amygdala, but not pre-SMA and hippocampus, was significantly higher in trials that were later correctly remembered (Fig. 5f). Together, this shows that the MFC contains neurons with persistent activity during maintenance. This persistent activity was not indicative of memory content but correlated with memory load and later performance.

Probe neurons

The third group of cells we characterized were probe neurons. These neurons increased their firing rate only during the presentation of the probe stimulus, relative to encoding and maintenance (Fig. 6a,b shows an example; see methods). The response of probe neurons was not visually selective: only 5 (6.49%) neurons out of 77 showed image selectivity ($P=0.385$ against scrambled data). Also, probe neurons did not respond to the identical stimuli when presented during encoding ($P=0.0005$, permutation test, Fig. 6d). We identified probe neurons in all areas, but most prominently in pre-SMA ($\chi^2[1]=44.21$; $P=2e-11$; Fig. 6c), where 26.2% of neurons were probe neurons. Next, we tested whether the response of these neurons is related to visual input or movement initiation by computing the maximal firing rate after aligning to stimulus onset and to button press (RT). Aligning the response with stimulus onset resulted in significantly higher peak firing rates (mean for all probe neurons: $8.82 \text{ Hz} \pm 8.74$) compared to the peak firing rate when aligned to button press ($5.85 \text{ Hz} \pm 7.5$, permutation t-test: $t(76)=8.71$; $P=0.0005$, see Fig. 6a, b for an example). This suggests that these neurons responded to the probe itself rather than to movement initiation.

What was the relationship between the activity of probe neurons and behavior? A minority of the probe neurons recorded in pre-SMA (15.6% of probe neurons) showed a differential response as a function of the WM-based decision (in vs. out, Fig. 6e, counting spikes in a window of -800ms to 0ms relative to button press). Similarly, a small number of neurons differentiated the button (left or right) used to communicate the decision (Fig. 6f). However, the large majority of probe neurons (84%) signaled neither. Instead, these neurons showed a strong but indiscriminate increase in firing rate to all pictures shown as probes. We hypothesize that probe neurons signal a change of task phase, i.e. a switch from maintenance to retrieval of information held in WM (see discussion).

Separability and effects of epilepsy

Were the three neuronal classes we identified separate or were there neurons which qualified as multiple types? Note that concept and maintenance cells are observed preferentially in separate anatomical areas (MTL vs. MFC), supporting the argument that they are separate. Nevertheless, there was some overlap in the same areas (Supplementary Fig. 4a). However, further analysis of the correlations between effect sizes attributable to different factors shows that the three types we identified constitute three largely separate categories of neurons (Supplementary Fig. 4b–c). Lastly, we confirmed that the effects observed did not differ between neurons recorded within vs. outside of later resected tissue (Supplementary Fig. 4d).

Decoding of information held in WM

We next assessed how well neuronal activity during maintenance predicted the identity of remembered images and/or performance in a given trial. For this, we trained a decoder on a pseudo-population of all neurons recorded of a given kind (for example, only MFC neurons or only concept cells) on a subset of trials and tested its performance on an independent set of test trials (see methods). We used this decoding approach to assess whether activity during maintenance was sufficiently strong to be read out at the single-trial level. We found that neurons in the MTL, but not MFC, carried information about picture identity during load 1 trials (Fig. 7a, green bars). Because there were multiple correct answers for loads 2 and 3, we used a separate binary decoder (trained in load 1 trials) for each image to predict whether it is currently in memory or not in load 2/3 trials, and averaged performance across all 5 images. This allowed us to compare performance between all load conditions. The average decoding accuracy for amygdala and hippocampal neurons during load 1 was 83.5 % and 63.4%, respectively (Fig. 7a, $P=0.002$ and $P=0.007$, respectively, estimated using scrambled data, see methods, chance was 50%). We observed a sharp drop in decoding accuracy as a function of increasing load (Fig. 7a) in both amygdala and hippocampus. Nevertheless, decoding performance was above chance for loads 2 and 3 in the amygdala (load 2: 60.7 %, $P=0.008$; load 3: 56 %, $P=0.039$) but not in the hippocampus (load 2: 54.2 %, $P=0.25$; load 3: 47.7 %, $P=0.73$). We also performed the same analysis using concept cells only (for load 1 trials). We found little difference in decoding accuracy between using only concept cells and the whole population of all recorded neurons (Fig. 7a, P for amygdala comparison = 0.88 and hippocampus = 0.32). Thus, the information decodable about picture identity from the population was carried principally by concept cells.

We next tested if information present in the neuronal firing rate of concept cells is carried by a static or dynamic code^{27,28}. We tested this by evaluating whether a decoder trained at one point of time could decode information obtained from a different point of time. We found that regardless of the time the decoder was trained, decoding performance was significantly above chance (Fig. 7b). For example, the average decoding accuracy during the maintenance period of a decoder trained based on the activity during the encoding period time (300 ms after image exposition) was 0.58 (chance level 0.2, $P=0.002$ vs. chance as estimated using a permutation test). Together, this shows that concept cells coded for the identity of the items held in memory using a static code. We next tested whether activity during maintenance was indicative of the number of items (load) held in memory. We found that decoding memory load was possible in both pre-SMA (decoding accuracy: 67.96 %, $P=0.002$, chance 33%, Fig. 7c) and dACC (42.88 %, $P=0.027$), but not Amygdala and hippocampus (Fig. 7c). When we restricted the decoder to only use information provided by the previously identified maintenance neurons, we found that load decoding performance in pre-SMA did not change significantly (59.8 % vs 67.96 %, $P=0.061$). In contrast, in dACC, allowing access only to maintenance neurons reduced decoding performance to chance levels (performance 38.04%, $P=0.164$). Thus, activity of MFC cells did not carry information about the content of the memory itself but rather about other aspects of the memory: the current load. Additionally, we observed that in pre-SMA maintenance neurons carried most of the information about memory load. Note that the activity of concept cells in the MTL

also varied as a function of load (Fig. 4c), but only for the preferred stimulus of a cell. In contrast, the activity of MFC cells varied as a function of load regardless of the stimulus held in memory (as the above decoding results demonstrate).

Finally, we tested if it was possible to predict the amount of time (referred to as RT) taken to make the in/out decision at the end of the trial based on the activity of neurons during the maintenance period (Fig. 7d). Stronger memories generally lead to faster and more accurate and confident decisions²⁹. We found evidence for this behavior in our task also, because incorrect decisions took longer (see behavioral results). Based on this rationale, we split correct trials into two groups based on the median RT of each patient (slow/fast groups). We did this separately for each load and in/out group (6 groups total) to account for the systematic RT differences that result from different loads and in/out decisions. We found that RT could be predicted from neurons located in pre-SMA (65.28 %, $P=0.002$, chance 50%, Fig. 7d) and amygdala (60.86 % $P=0.0027$, Fig. 7d) but not in dACC and hippocampus. Repeating the same analysis using only maintenance neurons revealed similar performance in pre-SMA (59.8%, $P=0.0157$ vs. chance). In amygdala, decoding based on maintenance neurons only reduced performance to chance levels (52.65 %, $P=0.2$).

Together, this shows that information about the content of working memory is provided by the firing rate of concept cells in the MTL (Fig. 7a) and that information about the quality of memories (load, later RT) is encoded predominantly by neurons in the MFC (Fig. 7c,d).

Working memory as attractors in state space

We next quantified the dynamics of neuronal activity during memory maintenance. An attractor is a location in state space, which in our case is formed by the activity of all recorded neurons. Here, an attractor is thus a pattern of firing rates. Theoretical work suggests that during WM maintenance the neural trajectories reside close to a point in this state space (the attractor) to maintain information in WM. This space could potentially assume a variety of patterns, including a static code²², dynamic patterns²⁸ or combination of both³⁰. Our goal was to determine whether the neuronal trajectories during the maintenance period were compatible with an attractor formed by a static code. We used demixed PCA (dPCA) as a dimensionality reduction technique³¹ to derive basis functions into which all neural activity was projected. The projection matrix of dPCA was computed only based on activity during encoding period 1 (first image) on all neurons. Afterwards, we projected the activity during the other encoding and the maintenance periods into this same space without recomputing the basis functions, making this analysis unbiased. We used dPCA with picture identity as the marginalized variable³¹. We found that a four dimensional space formed by four dPCs (id 1,2,3,5) separated the neural trajectories well (30.7% explained variance) between the five different image identities and had the highest percentage of variance attributed to picture identity (Fig. 8a–b, Supplementary Fig. 5, Supplementary Fig. 6d, see also supplementary video 1).

We first analyzed the rate of change (velocity) in this four dimensional neural state. As expected, velocity was highest during encoding due to the strong visual-onset transients of concept cells (Fig. 8c). In contrast, velocity was substantially reduced at baseline (fixation cross). Surprisingly, the amount of change in the activity during maintenance was reduced:

the velocity during maintenance was not significantly different from that during baseline (Fig. 8c, $P=0.233$). We next quantified the pairwise distance between the neural trajectories associated with different images held in memory (Fig. 8d; only for load 1 because pairwise distance is ill defined if several items are held in memory. See analysis using DA below for other loads). Distances were maximal during encoding but were maintained significantly above baseline levels during maintenance (Fig. 8d, $P=0.0005$). Thus, the neural trajectories during the maintenance period were drawn to particular locations in state space, which is the definition of an attractor^{21,32}. Therefore, our data suggests that neural activity during maintenance clusters around attractors (Fig. 8d).

We next tested whether the distance of the neural trajectories from a given attractor was correlated with behavior (accuracy and RT). We defined the attractor location based on the activity during maintenance of one item and quantified the distance of the neural activity to this activity in each trial using a distance metric (DA, “distance to attractor”, see methods). DA was defined as the Euclidian distance across trials between the attractor center for a particular image divided by the average distance of this attractor from all other attractors representing all the other pictures (Fig. 8e shows an illustration of this concept in 2d space). We found that the DA was smaller for trials in which the item in memory was correctly remembered compared to trials where the item was forgotten (Fig. 8f, permutation test: $P=0.002$, computed for loads 1–3 separately and averaged; see Supplementary Fig. 6c for each load). This difference was abolished when we excluded concept cells from the population, but excluding maintenance neurons did not alter this effect (Fig. 8f). Similarly, excluding neurons recorded in MFC did not affect this result whereas excluding MTL neurons eliminated the observed differences (Fig. 8f, permutation test, $P=0.002$ and $P=0.664$, respectively). We also tested whether the distance from the attractor predicted later decision time. We tested this for the subset of trials in which the probe was held in memory (“IN trials”). Faster decisions were associated with smaller DA values, indicating that neural trajectories (during maintenance) were closer to the given attractor (Fig. 8g). Similar to before, that was true only when concept cells and MTL were present in the population.

We next repeated the same analysis only using simultaneously recorded neurons ($n=5$ subjects who had at least four simultaneously recorded concept cells). Similar to the pseudo-population results, trial-by-trial, the distance to the closest attractor was predictive of whether the answer given by the subject in a trial would be correct or incorrect (Fig. 8h, permuted t-test, $t[624]=3.33$, $P=0.003$). This shows that attractors can be seen in individual subjects and that the distance to the attractor is informative in individual trials.

In conclusion, this shows that persistent activity during maintenance created attractors and that the distance of the neural trajectory in a given trial predicted the quality of the memory trace as measured by reaction time and accuracy.

Discussion

We found neurons with two different types of persistent activity in MTL and MFC during WM maintenance. In the MTL, concept cells were persistently active, and this activity was characterized by a high degree of selectivity. In the MFC, on the other hand, we observed a

group of neurons (Fig. 5c) which tonically increased their activity during maintenance, but this activity was indifferent to the identity of the memoranda currently held in the memory.

Preceding the WM task, we used a screening task to identify visually selective ‘concept cells’ similar to those first identified by Quiroga et al^{24,33–35}. Such concept cells responded strongly and selectively only to the preferred, but not to the non-preferred images (see Supplementary Fig. 2a). Interestingly, we identified more concept cells in the amygdala (25.4%) compared to the hippocampus (12.1%), a trend compatible with that observed previously³⁴. It remains an open question whether concept cells also exhibit persistent activity for stimuli other than images, such as text or audio³⁴. In addition to transient activation by visual input, concept cells can also be activated by free recall of episodic memories³⁶ and visual imagery³⁷. While this has so far been interpreted as representing recall from long-term memory, we here now show a direct relationship between persistent concept cell activity and WM. Of note, we found that the response strength (firing rate) of concept cells for images held in memory decreased as a function of memory load, an observation that fits models of persistent activity³⁸.

We found that the dynamics of self-sustained persistent activity were well described as an attractor state. Attractor-like behavior has been observed in animal experiments as well as theoretical work^{20–22,39}. Modeling work shows that attractor networks can hold multiple pieces of information encoded by separate attractors⁴⁰ in a robust and noise-resistant manner³⁹. Theoretical work further indicates that the activity of attractor networks can drift during the maintenance period due to other attractors located close-by in state space^{20,39}. This prediction has been confirmed experimentally in macaque recordings with a single item held in memory, where the extent of drift away from an attractor was predictive of errors²². Here, we present evidence that the activity of human neurons during WM maintenance is compatible with that expected if attractors were present. We show that this prediction of the model holds for discrete attractors and for multiple items held in memory, because we find that the quality of a memory trace can be predicted by the drift (distance) away from the attractor representing a given image (Fig. 8). In addition, we found that it is principally concept cells in the MTL that support attractors (Fig. 8 f , g). Here, we defined each attractor based on the pattern of activity observed when only a single item was held in memory. However, when multiple items were held in memory (loads 2–3), multiple concept cells are persistently active. Consequently, the population activity is different from that defined by the single-item attractors as demonstrated by increased DA values (Supplementary Fig. 6c). Despite this, the distance to the attractors defined in this manner was predictive of memory content and quality even in loads 2 and 3 (Supplementary Fig. 6c). This is possible because of the high dimensionality of the state space, which make it possible to be closest to several attractors while being far away from all the others⁴⁰. Overall, our findings reveal a direct neuronal correlate for discrete attractors in humans. This shows that persistent activity can be viewed as forming attractors and that persistent activity is a mechanism capable of supporting information maintenance during WM.

A lack of evidence for stimulus-specific persistent activity in some studies^{28,41} has led to alternative proposals on how WM is maintained. A key new proposal is that information in WM can be carried by salient states encoded through synaptic weight changes^{42,43} or by

oscillatory bursts⁴¹. Our analysis now shows that persistent activity compatible with the original model exists in the human MTL. This argues that the original WM model is a valid way to conceptualize persistent activity during Sternberg-like WM tasks. However, this does not rule out the possibility that mechanisms other than persistent activity also contribute to WM.

Our data shows that the activity of MTL neurons during WM maintenance carried information about the items held in memory. This raises the question of the role of the MTL for WM. While it has been observed before that intracranial field potentials recorded from the MTL are modulated by load^{18,44}, these differences are not stimulus selective. Subjects with bilateral MTL damage do not exhibit WM deficits in some circumstances¹⁷, but sufficiently difficult WM tasks or the presence of distractors¹⁷ reveal WM deficits in the same subjects. Based on this, it has been suggested that the MTL is critical in conditions when WM alone is insufficient¹⁷. A possible role for the persistent MTL activity we observed is therefore to maintain memory engrams so that they can be used to recover information in WM that has been lost due to shifts in attention triggered by distractors or when the capacity of WM is saturated.

We found amygdala neurons with prominent persistent activity. Interestingly, amygdala persistent activity was more easily decodable than activity in the hippocampus. This is in contrast to episodic memory-related activity during free recall³⁶, which is prominent in the hippocampus but not the amygdala. This raises the novel possibility that the amygdala could have a role in supporting WM. While encoding new long-term memories does not require the amygdala, the amygdala prominently modulates this process by enhancing or suppressing the encoding of new memories⁴⁵. Interestingly, patients with amygdala lesions frequently have working memory deficits⁴⁶. Also, rats with amygdala lesions have a specific impairment in visual WM tasks⁴⁷. While it is recognized that the amygdala can boost general vigilance and thereby facilitates execution of demanding tasks⁴⁵, our data indicates that the amygdala might directly contribute to WM maintenance by storing specific memoranda with the help of persistently active neurons.

The persistent activity of MFC neurons did not carry information about the content of the memory. Nevertheless, this activity was relevant for WM, because it was predictive of later behavior (accuracy and RT) and memory load. This is similar to some recordings in macaques, which also revealed non-selective persistent activity in the MFC¹². However, other macaque studies have revealed MFC single-units who are stimulus selective^{11,48}. It remains an open question whether stimulus-specific persistent activity in the human MFC exists in other tasks or for other stimuli such as tactile vibration stimuli¹¹, sounds¹¹ or task contexts⁴⁸. Note that we did not record from dlPFC in this study, which is frequently associated with stimulus selective persistent activity^{4,5,7}. It therefore remains an open question whether human dlPFC neurons have stimulus-specific persistent activity. Nevertheless our finding that the activity of maintenance neurons in the MFC during delay was predictive of RT and accuracy is consistent with the hypothesized role of the MFC in attentional control⁴⁹.

dACC and pre-SMA are both part of the cingulo-opercular (CO) system. It is thought that the activity of this network supports implementation of different task sets¹⁵. Indeed, in lesion studies, subjects with MFC lesions are impaired in switching between different response and instruction sets⁵⁰. Our task consisted of three task sets (encoding, maintenance, retrieval). A role of the MFC in controlling switches between these three task sets is compatible with our finding that differences in activity due to load and RT appear prominently only in the initial part of the maintenance period (Fig. 5a, b). This is because in each trial, the set size is unknown to the subject. Therefore, starting the maintenance period after encoding only one picture was more unpredictable compared to starting it after encoding three pictures. Thus, switching to maintenance after encoding one picture would be more attentionally demanding. Indeed, we found that pre-SMA cells were most active for load 1, and less so for loads 2 and 3 (Fig. 5). Similarly, the firing rate of MFC maintenance neurons was most predictive of RT early during the maintenance period (Fig. 5b), which might also be related to switching from the encoding to the maintenance mode. Here, we suggest that these neurons support attentional control and, together with probe neurons (see below), the implementation of task sets. Future work is needed to determine whether maintenance- and concept neurons interact during task.

Finally, we also observed a group of “probe neurons” in pre-SMA that only became active during the retrieval phase of the task. The activity of these neurons was better explained by the time of visual stimulus onset rather than the timing or nature of the motor response. Probe neurons did not respond to the same stimuli during encoding, demonstrating that they are not simply visually tuned or responsive. In the context of task sets¹⁵, “probe neurons” signaled the transition to the retrieval phase of the task, i.e. a switch from maintenance to retrieval of information held in WM. Together, this reveals a direct correlate of neurons implementing transitions between different tasks sets in human MFC.

Online Methods

Task

We used a modified Sternberg task with images (instead of the usual digits) as material for memorization (Fig. 1a). Each trial started with a fixation cross shown for 900 – 1000 ms. Next, we sequentially presented the images to be memorized (“encoding”) in a given trial. Each picture was presented for 1 second, followed by a blank screen for 1–200ms (randomized). Subjects were asked to memorize the 1–3 images shown in each trial. We use the terms “encoding 1”, “encoding 2” and “encoding 3” to refer to the 1–3 images shown in a trial. After encoding, there was a maintenance (delay) period lasting at least 2.5s and at most 2.8s. During this time, the word “hold” was shown on the screen. Lastly, after the end of the maintenance period, a probe stimulus was displayed. Subjects were asked to decide if the probe stimulus was shown as one of the preceding 1–3 images or not. Participants responded by pressing the “green” or “red” buttons on a response pad. Which color was “yes” and which “no” was shown at the top of the screen during each probe trial. We used this approach to switch the location of the yes/no buttons in the middle of the experiment as a control. We asked subjects to respond as fast as possible. The probe picture was presented

until subjects made a response. In each session, subjects preformed 108 or 135 trials depending on the task variant. Pictures were shown in pseudo-random order.

The pictures used for each participant were different and were determined by the results of a screening task conducted 2–3 hours earlier (except from one patient who did not undergo a screening procedure and had images used with a previous subject). For the screening task, images were chosen based on a patient's interests.

During screening, we showed 54–64 images. Each image was shown 6 times in randomized order for 1s. As a control, every few trials (randomized) we asked a control question related to the image shown immediately before (i.e. "Did the last image present person/ landscape/ animal?"). After the screening task, we immediately analyzed the data to choose the 5 images with the best responses as judged by the mean response in a 200–1000ms window relative to stimulus onset (based on an F statistic computed by a one-way ANOVA with image as a factor). In sessions in which less than 5 neurons showed significantly selective responses (6/13 sessions), we picked the remaining images according to the strongest non-selective responses. In these instances, we used a combination of the F-statistic (with image as a factor) and neuronal isolation quality (amplitude of waveform) to choose the next best non-selective response. These 5 images were subsequently used for the Sternberg task. Both tasks were implemented in MATLAB using the Psychophysics Toolbox ⁵¹. Note that the statistical tests run during the screening and Sternberg task were statistically independent, because they were run at distinct periods of time.

Data collection and analysis were not performed blind to the conditions of the experiments.

Patients

13 subjects participated in the study (Supplementary Table 1). All of them were implanted with depth electrodes for possible surgical treatment of epilepsy. They volunteered for the study and gave informed consent. This study was approved by the Institutional Review Boards of the Cedars-Sinai Medical Center, Huntington Memorial Hospital, and the California Institute of Technology. Electrodes were localized based on pre-and post-operative T1 structural MRIs. We used the following processing pipeline to transform the post-operative MRI into the same space as a template brain. We extracted the brains from the pre-and post-operative T1 scans ⁵² and aligned the post-operative to the pre-operative scan with Freesurfer's `mri_robust_register` ⁵³. We then computed a forward mapping of the pre-operative scan to the CIT168 template brain ⁵⁴ using a concatenation of an affine transformation followed by a symmetric image normalization (SyN) diffeomorphic transform computed by the ANTs software package ⁵⁵. This resulted in a post-operative scan overlaid on the MNI152-registered version of the CIT168 template brain ⁵⁴. We then used Freesurfer's Freeview program to mark the electrodes as point sets to determine where the tips of the microwires were located. For visualization only, electrode locations were projected onto the 2D sagittal plane (Fig. 1b).

Spike sorting and quality metrics of single units

Each macroelectrode contained eight 40 μm diameter microwires⁵⁶. We recorded broadband (0.1–9000Hz filter) from a total of 64 channels sampled at 32 kHz using a Neuralynx Atlas system. Signals were locally referenced to one of the eight microwires in each brain area.

The raw signal was filtered with a zero-phase lag filter in the 300–3000Hz band and spikes were detected and sorted using a semi-automated template-matching algorithm²³. We computed several spike sorting quality metrics for all identified putative single units to assess the quality of identified units (Supplementary Fig. 1): 1) the percentage of interspike intervals (ISIs) below 3 ms was $0.42\% \pm 0.68\%$, 2) the ratio between the standard deviation of the noise and the peak amplitude of the mean waveform of each cluster was 6.18 ± 4.09 (peak SNR), 3) the pairwise projection distance in clustering space between all neurons isolated on the same wire was 14.39 ± 6.2 (projection test; in units of s.d.⁵⁷ of the signal), 4) the modified coefficient of variation of variability in the ISI (CV2) was 0.94 ± 0.14 , 5) the median isolation distance⁵⁸ was 30.9. We calculated the isolation distance in a ten-dimensional feature space (energy, peak amplitude, total area under the waveform and first five principal components of the energy normalizes waveforms⁵⁸).

Statistics

Statistical comparisons were conducted using permutation tests based on a null distribution estimated from $B=2000$ runs on data with scrambled labels using the EEGLAB toolbox⁵⁹. For comparisons with two groups we used the permuted t-test statistic and for comparisons with more than two groups the estimated distribution was that of the F statistics. Note that we used permutation tests throughout to avoid the assumption of normality. We used ANOVAs instead of linear regressions to test for effects of load because linear regressions introduce the additional assumption of a monotonic relationship between load and behavior. Note that the reported p-values can be different from those expected from the t and F distributions because p-values were based on the empirically estimated null distribution. No statistical methods were used to pre-determine sample sizes but our sample sizes are similar to those reported in previous publications in the field⁵⁶.

Selection of neurons

For each recorded neuron, we ran the following three statistical tests to determine whether a cell qualified as a concept, maintenance, or probe neuron. Tests were independent and some cells qualified as multiple types. To identify concept cells we counted spikes in a window 200–1000ms following stimulus onset of the first encoding period. Concept cells were selected ($p < 0.05$) using a permuted one-way ANOVA with 5 groups (the number of unique images). In addition, we also required that the activity for the image with the maximal response was significantly larger compared to that of all other images ($p < 0.05$, permutation t-test). The value of this latter test was also used as the effect size metric to compare the neuronal groups (see “Separability of neuronal categories”). For the analysis of persistent activity of concept cells, we used all concept cells found during the Sternberg task regardless of whether they were also significant in the screening task.

To identify maintenance cells, we tested if the mean firing rate during the maintenance period (0–2500ms) was significantly larger (permutation t-test, $p < 0.05$) than the firing rate during the 500ms long fixation cross period before trial start. If a maintenance cell was also a concept cell, we in addition required that its maintenance activity for the non-preferred stimuli was significantly higher than baseline. This second criteria is to assure that strong concept cells, which only increase their firing rate to a single stimulus, don't automatically qualify as maintenance cells.

Probe neurons were identified by comparing the firing rate during presentation of the probe stimulus (spikes counted during 200–1000 ms after probe onset) to that during both the encoding ($p < 0.05$, 200 – 1000 ms) and maintenance periods ($p < 0.05$, 0–2500 ms, permutation t-test).

Single-cell metrics

We used a Picture Selectivity Index (PSI) to quantify the response of concept cells:

$$\text{PSI} = \frac{\text{preferred image FR} - \text{not preferred image FR}}{\text{Baseline FR}} * 100$$

Where the FR is the mean firing rate in a 200–1000 ms window relative to stimulus onset for encoding and probe, 0–2500 ms for maintenance and –500-0 ms for baseline. When comparing correct and incorrect trials, we computed firing rates and PSI values independently for the three load conditions (1,2,3) and then averaged the results in order to eliminate bias due to different number of incorrect trials for each load. We also used the depth of selectivity index ⁶ to quantify the sparsity of concept cells (Supplementary Fig. 2):

$$S = \frac{n - \left(\frac{\sum r_i}{r_{\max}} \right)}{n - 1}$$

Where n is the number of images presented, R_i is the firing rate of the neuron during the presentation of the i th picture, and R_{\max} is the largest firing rate across all presented images. A neuron with a S value of 0 would respond identically to all images, whereas a neuron with $S=1$ would respond only to a single image and not at all to others.

In order to compare time courses of neuronal activity, we counted spikes in 200 ms bins moved by a step-size of 2ms. We corrected for multiple comparisons using a cluster based approach ⁶⁰. In this method, a cluster is defined as a group of adjacent significant tests. We tested, for all identified clusters, if the summed value of the test statistic in a given cluster was larger than the 95th percentile of the same value estimated from scrambled labels. For comparing firing rates across cells, we standardized the firing rate of each neuron using the mean and standard deviation of the firing rate of each cell during the baseline (fixation cross, –500ms – 0 ms window before encoding 1). To evaluate if the proportion of significant neurons in an area was larger than that expected by chance we computed a null distribution based on randomly scrambled labels ($B=500$) and then estimated an empirical P value as described.

Population decoding

We used a pseudo-population of neurons pooled across all recording sessions and subjects to determine decoding performance. To pool images from different subjects we arbitrarily labeled the individual images of each subjects as image A-E. We then pooled the images with the same label across all subjects. Note that each subject saw a different set of pictures. Nevertheless, neurons can be considered a pseudo-population for decoding purposes, a procedure that is frequently used to pool single-neuron recordings across sessions and animals ⁶¹. This is because for the learning algorithm, all that matters is that each patient saw five distinct images regardless of their identity. The decoder assigns a weight to each neuron independently, i.e. the decoder will identify all neurons which signaled the presence of “image A” in all patients, regardless of the identity of “image A”. This pooling procedure provides, at every point of time, a $N \times T$ matrix where N is the total number of neurons recorded from all subjects in given brain area, and T is a smallest number of trials of a given type observed in all patients. This matrix thus represented, as a function of time, the neural state in an N dimensional space. We then trained a decoder to separate patterns in this high dimensional space. Spikes were counted in 1500 ms bins moved with a stepsize of 250 ms. For decoding, we used a support vector machine (SVM) as implemented in the `ndt` toolbox ⁶¹ and the `LIBSVM` library ⁶². We used leave-one-out cross-validation to estimate performance: one trial was randomly assigned for each neuron from each class as the test trial and the remaining trials were used for training. All possible train/test combinations were computed, and we used 500 randomly chosen train/test combinations to estimate the cross-validated testing error. For decoding image identity, we trained a classifier on the load 1 condition and then applied it to loads 2 and 3 to determine the accuracy by which the algorithm could determine whether a given image was in held in memory. To test if given decoding performance was significantly better than chance, we created null distributions by decoding using the same approach as described above but with scrambled labels. We repeated this computation 500 times. Note that the minimal possible P value is thus 0.002 and we used this p -value if no value in the null distribution exceeded the observed value. Similarly, for comparisons between decoding accuracy of different decoders - for example, all neurons compared to maintenance only neurons - we created a null distribution of differences in a given comparison by subtracting the null distribution created for the whole population (in this example) from the null distribution created for maintenance neurons only. We then estimated the p value for this comparison by counting the number of times the null distribution exceeded the observed difference of performance between the two decoders.

Analysis of neuronal activity dynamics

To analyze the dynamics of activity of the entire neuronal population we reduced dimensionally using demixed principal component analysis (dPCA) ²⁹. We used dPCA instead of PCA because PCA stretches components “blindly” based on the percentage of explained variance as a criterion. Due to this, the dimensions chosen often don’t have a meaningful interpretation and thus don’t help in a particular analysis. In contrast, dPCA has the advantage that it stretches components along dimensions not only to explain overall signal variance but also to explain variance attributable to variables of interest (such as image identity in our case). For dPCA analysis, we used a pseudo-population of all recorded neurons (as defined above). We used dPCA with picture identity and time as the

marginalized variable. We binned neuronal firing in 2ms non-overlapping bins and smoothed the resulting time course by convolution with a Gaussian kernel (200 ms width). In addition, we z-scored all time-courses based on the mean and standard deviation estimated from the baseline (-500-0ms relative to onset of first image). We computed the basis functions of dPCA (demixing weights) based only on the data recorded during 200-1000ms following onset of the first image ("encoding 1"). We rank-ordered the demixed principal components (dPCs) by their explained variance and used the first 15 dPCs. These together accounted for 51.58 % of the total signal variance. To prevent overfitting, we used a regularization procedure to find the optimal lambda parameter (Supplementary Fig. 6f). In addition, we also tested if the percentage of explained variance for picture identity was higher than that observed for data with scrambled image identity labels. Indeed, the portion of variance that was attributed to image identity by the 15 largest dPCs of the real data was 45.13%, whereas for the scrambled data it was only 13.64% (permutation test, $P=0.002$, Supplementary Fig. 6e). We projected the data from the maintenance period onto the basis functions computed from the encoding 1 period.

We used the multidimensional Euclidean distance $d(p,q)$ to quantify how different the population activity was between the neural activity vector $p(t)$ and $q(t)$ (which here are the neuronal states during two different conditions at time t):

$$d(p(t), q(t)) = \sqrt{\sum_{i=1}^n (p_i(t) - q_i(t))^2}$$

In addition, we quantified the speed by which the population activity changed at a given point of time t as $V(t)$:

$$V(t) = \frac{1}{n} * \frac{\sqrt{\sum_{i=1}^n (p_i(t) - p_i(t - \Delta t))^2}}{\Delta t}$$

Where n is the number of dimensions in dPCA space ($n=4$), t is 50 ms and p is the neural activity vector. We used the first four dPCs (id 1,2,3,5) that had the highest percentage of variance attributed to picture identity during load 1 trials (Fig. 8a-b, Supplementary Fig. 5, Supplementary Fig. 6d, see also supplementary video 1). dPC 4, in contrast, had variance that was only attributable to time but not image identity (Supplementary Fig. 5). Also, the next biggest dPCs (id 6) accounted for only 2.55% of the variance attributed to picture identity. We also recomputed our analysis using the 8 or 12 dPCs with the highest percentage attributed to image identity and found that results were very similar to the analysis based on 4 dPCs (Supplementary Fig. 6ab). Both V and d are population-level metrics not computable for single neurons. To estimate the variance of V and d , we bootstrapped the confidence intervals of V and d by randomly picking a subset 10% of trials and computed V and d for each such subset (repeated 50 times).

For each possible remembered stimulus k , we defined the location of its corresponding attractor A_k as the center (mean across time) of the neuronal trajectory observed during

maintenance of that image during the load 1 condition. There were 5 attractors (a–e) — one for each image used. To quantify the distance of a neuronal trajectory from an attractor in a given condition, we used a “distance to attractor” (DA) metric:

$$DA(t) = \frac{d_i(t)}{\frac{1}{C-1} \sum_{j \neq i}^C d_j(t)}$$

Where $d_j(t) = d(A_j, p(t))$ is the Euclidean distance of the neuronal trajectory $p(t)$ from the attractor A_j . C is the total number of attractors ($C=5$). We computed $DA(t)$ separately for loads 1–3 and each time point t during maintenance (0–2500 ms) and then averaged $DA(t)$ over all points of time and loads to get a single value DA for every trial. Note that $DA=1$ implies that the neural trajectory was equidistant between the tested (d_i) and all other attractors, which indicates no memory. On the other hand, $DA < 1$ indicates that the neuronal state was closer to one attractor compared to all the other attractors.

Similarly to d and V , DA is a population-level metric not computable for single neurons. We thus assessed the significance of differences in DA between correct and incorrect trials and fast and slow RT conditions by estimating a null distribution of DA based on data with scrambled labels. For visualization of distances between attractors in 2D space (Fig. 8e), we used non-metric multidimensional scaling (MDS). We used Euclidean distance as the pairwise distance measure. MDS was used for visualization only.

For the single subject analysis (Fig. 8h), we used the same approach as described for the pseudo-population but used only neurons recorded simultaneously from a given subject. For each session, we used the first four dPCs that explained the highest proportion of variance and computed DA accordingly.

Data availability

The data that support the findings of this study are available on reasonable request from the corresponding author. The data are not publicly available because they contain information that could compromise research participant privacy/consent.

Code availability

Analysis was performed in Matlab using the publicly available software packages *Osort*, *ndt*, and *EEGLAB* together with custom-developed analysis routines.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

We thank J. Minxa, R. Adolphs and J. Dubois for discussion, and the staff and physicians of the Epilepsy Monitoring Unit at Cedars-Sinai Medical Center and the Huntington Memorial Hospital for invaluable assistance. This work was supported by the National Science Foundation (1554105 to U.R.), the National Institute of Mental Health (R01MH110831 to U.R.), the McKnight Endowment Fund for Neuroscience (to U.R.), a NARSAD Young Investigator grant from the Brain & Behavior Research Foundation (23502 to U.R.), and the Pfeiffer Foundation (to U.R.).

References

1. Baddeley A. Working memory: theories, models, and controversies. *Annu. Rev. Psychol.* 2012; 63:1–29. [PubMed: 21961947]
2. Eriksson J, Vogel EK, Lansner A, Bergström F, Nyberg L. Neurocognitive Architecture of Working Memory. *Neuron.* 2015; 88:33–46. [PubMed: 26447571]
3. Goldman-Rakic PS. Cellular basis of working memory. *Neuron.* 1995; 14:477–485. [PubMed: 7695894]
4. Constantinidis C, Franowicz MN, Goldman-Rakic PS. The sensory nature of mnemonic representation in the primate prefrontal cortex. *Nat. Neurosci.* 2001; 4:311–316. [PubMed: 11224549]
5. Funahashi S, Bruce CJ, Goldman-Rakic PS. Mnemonic coding of visual space in the monkey's dorsolateral prefrontal cortex. *J. Neurophysiol.* 1989; 61:331–349. [PubMed: 2918358]
6. Rainer G, Asaad WF, Miller EK. Selective representation of relevant information by neurons in the primate prefrontal cortex. *Nature.* 1998; 393:577–579. [PubMed: 9634233]
7. Chafee MV, Goldman-Rakic PS. Matching patterns of activity in primate prefrontal area 8a and parietal area 7ip neurons during a spatial working memory task. *J. Neurophysiol.* 1998; 79:2919–2940. [PubMed: 9636098]
8. Fuster JM, Jervey JP. Inferotemporal neurons distinguish and retain behaviorally relevant features of visual stimuli. *Science.* 1981; 212:952–955. [PubMed: 7233192]
9. Miller EK, Li L, Desimone R. Activity of neurons in anterior inferior temporal cortex during a short-term memory task. *J. Neurosci.* 1993; 13:1460–1478. [PubMed: 8463829]
10. Suzuki WA, Miller EK, Desimone R. Object and place memory in the macaque entorhinal cortex. *J. Neurophysiol.* 1997; 78:1062–1081. [PubMed: 9307135]
11. Vergara J, Rivera N, Rossi-Pool R, Romo R. A Neural Parametric Code for Storing Information of More than One Sensory Modality in Working Memory. *Neuron.* 2016; 89:54–62. [PubMed: 26711117]
12. Constantinidis C, Procyk E. The primate working memory networks. *Cogn. Affect. Behav. Neurosci.* 2004; 4:444–465. [PubMed: 15849890]
13. Fuster JM. Cortex and memory: emergence of a new paradigm. *J. Cogn. Neurosci.* 2009; 21:2047–2072. [PubMed: 19485699]
14. Yarkoni T, Poldrack RA, Nichols TE, Van Essen DC, Wager TD. Large-scale automated synthesis of human functional neuroimaging data. *Nat. Methods.* 2011; 8:665–670. [PubMed: 21706013]
15. Dosenbach NUF, et al. A Core System for the Implementation of Task Sets. *Neuron.* 2006; 50:799–812. [PubMed: 16731517]
16. Squire LR, Stark CEL, Clark RE. The medial temporal lobe. *Annu. Rev. Neurosci.* 2004; 27:279–306. [PubMed: 15217334]
17. Jenson A, Squire LR. Working memory, long-term memory, and medial temporal lobe function. *Learn. Mem.* 2012; 19:15–25. [PubMed: 22180053]
18. Axmacher N, et al. Sustained neural activity patterns during working memory in the human medial temporal lobe. *J. Neurosci.* 2007; 27:7807–7816. [PubMed: 17634374]
19. Ojemann GA, Creutzfeldt O, Lettich E, Haglund MM. Neuronal activity in human lateral temporal cortex related to short-term verbal memory, naming and reading. *Brain.* 1988; 111:1383–1403. [PubMed: 3208062]
20. Wang XJ. Synaptic reverberation underlying mnemonic persistent activity. *Trends in Neurosciences.* 2001; 24:455–463. [PubMed: 11476885]
21. Hopfield JJ. Neural networks and physical systems with emergent collective computational abilities. *Proc. Natl. Acad. Sci. U. S. A.* 1982; 79:2554–2558. [PubMed: 6953413]
22. Wimmer K, Nykamp DQ, Constantinidis C, Compte A. Bump attractor dynamics in prefrontal cortex explains behavioral precision in spatial working memory. *Nat. Neurosci.* 2014; 17:431–439. [PubMed: 24487232]
23. Rutishauser U, Douglas R. State-dependent computation using coupled recurrent networks. *Neural Comput.* 2009; 509:478–509.

24. Quian Quiroga R, Reddy L, Kreiman G, Koch C, Fried I. Invariant visual representation by single neurons in the human brain. *Nature*. 2005; 435:1102–1107. [PubMed: 15973409]
25. Mormann F, et al. A category-specific response to animals in the right human amygdala. *Nat. Neurosci*. 2011; 14:1247–1249. [PubMed: 21874014]
26. Sternberg S. In defence of high-speed memory scanning. *Q. J. Exp. Psychol*. 2016; 69:2020–2075.
27. Meyers EM, Freedman DJ, Kreiman G, Miller EK, Poggio T. Dynamic population coding of category information in inferior temporal and prefrontal cortex. *J. Neurophysiol*. 2008; 100:1407–1419. [PubMed: 18562555]
28. Stokes MG, et al. Dynamic coding for cognitive control in prefrontal cortex. *Neuron*. 2013; 78:364–375. [PubMed: 23562541]
29. Rutishauser U, et al. Representation of retrieval confidence by single neurons in the human medial temporal lobe. *Nat. Neurosci*. 2015; 18:1–12. [PubMed: 25547471]
30. Murray J, et al. Stable population coding for working memory coexists with heterogeneous neural dynamics in prefrontal cortex. *Proc. Natl. Acad. Sci*. 2016
31. Kobak D, et al. Demixed principal component analysis of neural population data. *Elife*. 2016; 5:1–37.
32. Milnor J. On the concept of attractor. *Commun. Math. Phys*. 1985; 99:177–195.
33. Mormann F, et al. Latency and selectivity of single neurons indicate hierarchical processing in the human medial temporal lobe. *J. Neurosci*. 2008; 28:8865–8872. [PubMed: 18768680]
34. Quian Quiroga R, Kraskov A, Koch C, Fried I. Explicit encoding of multimodal percepts by single neurons in the human brain. *Curr. Biol*. 2009; 19:1308–1313. [PubMed: 19631538]
35. Waydo S, Kraskov A, Quian Quiroga R, Fried I, Koch C. Sparse representation in the human medial temporal lobe. *J. Neurosci*. 2006; 26:10232–10234. [PubMed: 17021178]
36. Gelbard-Sagiv H, Mukamel R, Harel M, Malach R, Fried I. Internally generated reactivation of single neurons in human hippocampus during free recall. *Science*. 2008; 322:96–101. [PubMed: 18772395]
37. Kreiman G, Koch C, Fried I. Imagery neurons in the human brain. *Nature*. 2000; 408:357–361. [PubMed: 11099042]
38. Macoveanu J, Klingberg T, Tegnér J. A biophysical model of multiple-item working memory: A computational and neuroimaging study. *Neuro Science*. 2006; 141:1611–1618.
39. Camperi M, Wang XJ. A model of visuospatial working memory in prefrontal cortex: Recurrent network and cellular bistability. *J. Comput. Neurosci*. 1998; 5:383–405. [PubMed: 9877021]
40. Laing CR, Troy WC, Gutkin B, Ermentrout GB. Multiple bumps in a neuronal model of working memory. *Siam J. Appl. Math*. 2002; 63:62–97.
41. Lundqvist M, et al. Gamma and Beta Bursts Underlie Working Memory. *Neuron*. 2015
42. Stokes MG. ‘Activity-silent’ working memory in prefrontal cortex: a dynamic coding framework. *Trends Cogn. Sci*. 2015; 19:394–405. [PubMed: 26051384]
43. Mongillo G, Barak O, Tsodyks M. Synaptic theory of working memory. *Science*. 2008; 319:1543–1546. [PubMed: 18339943]
44. van Vugt MK, Schulze-Bonhage A, Litt B, Brandt A, Kahana MJ. Hippocampal gamma oscillations increase with memory load. *J. Neurosci*. 2010; 30:2694–2699. [PubMed: 20164353]
45. Davis M, Whalen PJ. The amygdala: vigilance and emotion. *Mol. Psychiatry*. 2001; 6:13–34. [PubMed: 11244481]
46. Buchanan, TW., Tranel, D., Adolphs, R. The Human Amygdala. Whalen, PJ., Phelps, EA., editors. 2009. p. 289-317.
47. Peinado-Manzano MA. The role of the amygdala and the hippocampus in working memory for spatial and non-spatial information. *Behav. Brain Res*. 1990; 38:117–134. [PubMed: 2363833]
48. Saez A, Rigotti M, Ostojic S, Fusi S, Salzman CD. Abstract Context Representations in Primate Amygdala and Prefrontal Cortex. *Neuron*. 2015; 87:869–881. [PubMed: 26291167]
49. Bush G, Luu P, Posner MI. Cognitive and emotional influences in anterior cingulate cortex. *Trends Cogn Sci*. 2000; 4:215–222. [PubMed: 10827444]
50. Gläscher J, et al. Lesion mapping of cognitive control and value-based decision making in the prefrontal cortex. *Proc. Natl. Acad. Sci. U. S. A*. 2012; 109:14681–14686. [PubMed: 22908286]

Methods-only References

51. Brainard DH. The Psychophysics Toolbox. *Spat. Vis.* 1997; 10:433–436. [PubMed: 9176952]
52. Ségonne F, et al. A hybrid approach to the skull stripping problem in MRI. *Neuro image.* 2004; 22:1060–1075. [PubMed: 15219578]
53. Reuter M, Rosas HD, Fischl B. Highly accurate inverse consistent registration: A robust approach. *Neuro image.* 2010; 53:1181–1196. [PubMed: 20637289]
54. Tyszka JM, Pauli WM. A high resolution in vivo MRI atlas of the adult human amygdaloid complex. *Hum. Brain Mapp.* 2016; 37:3979–3998. [PubMed: 27354150]
55. Avants B, et al. Multivariate Analysis of Structural and Diffusion Imaging in Traumatic Brain Injury. *Acad. Radiol.* 2008; 15:1360–1375. [PubMed: 18995188]
56. Rutishauser U, Ross IB, Mamelak AN, Schuman EM. Human memory strength is predicted by theta-frequency phase-locking of single neurons. *Nature.* 2010; 464:903–907. [PubMed: 20336071]
57. Pouzat C, Mazor O, Laurent G. Using noise signature to optimize spike-sorting and to assess neuronal classification quality. *J. Neurosci. Methods.* 2002; 122:43–57. [PubMed: 12535763]
58. Harris KD, Henze DA, Csicsvari J, Hirase H, Buzsáki G. Accuracy of tetrode spike separation as determined by simultaneous intracellular and extracellular measurements. *J. Neurophysiol.* 2000; 84:401–414. [PubMed: 10899214]
59. Delorme A, Makeig S. EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *J Neurosci Methods.* 2004; 134:9–21. [PubMed: 15102499]
60. Maris E, Oostenveld R. Nonparametric statistical testing of EEG- and MEG-data. *J Neurosci Methods.* 2007; 164:177–190. [PubMed: 17517438]
61. Meyers EM. The Neural Decoding Toolbox. *Front. Neuroinform.* 2013; 7:8. [PubMed: 23734125]
62. Fan R-E, Chen P-H, Lin C-J. Working Set Selection Using Second Order Information for Training Support Vector Machines. *J. Mach. Learn. Res.* 2005; 6:1889–1918.

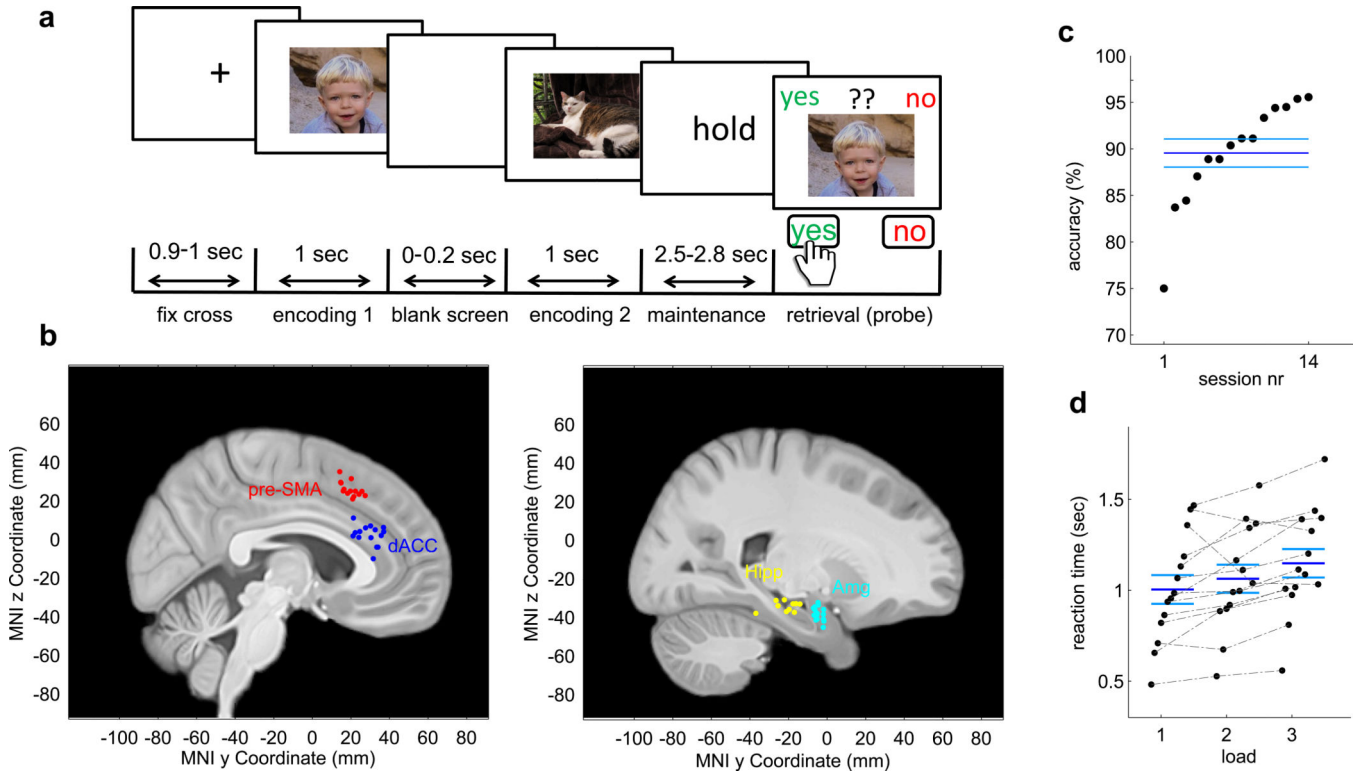


Figure 1. Task, recording locations, and behavioral results

(a) The task. Upper row represents examples of the screens presented to the subjects during an example trial (with load 2). The lower row represents the lengths of time for which each screen was shown. Each trial consisted of 1–3 sequentially presented pictures (encoding), followed by a variable delay (holding or maintenance period). After the delay, a probe image was shown and patients indicated whether the probe was or was not shown during the immediately preceding encoding period. (b) Location of recording sites in MNI152 space (see methods). Recording locations are indicated by different colors (red is pre-SMA, blue is dACC, yellow is hippocampus, and cyan is amygdala). (c–d) Behavioral results. (c) Accuracy of all sessions, rank-ordered. (d) Median reaction time (relative to onset of the probe image) as a function of load. Each dashed line connects an individual session. (b–c) Thick and light blue lines represent the mean and s.e.m across all sessions, respectively.

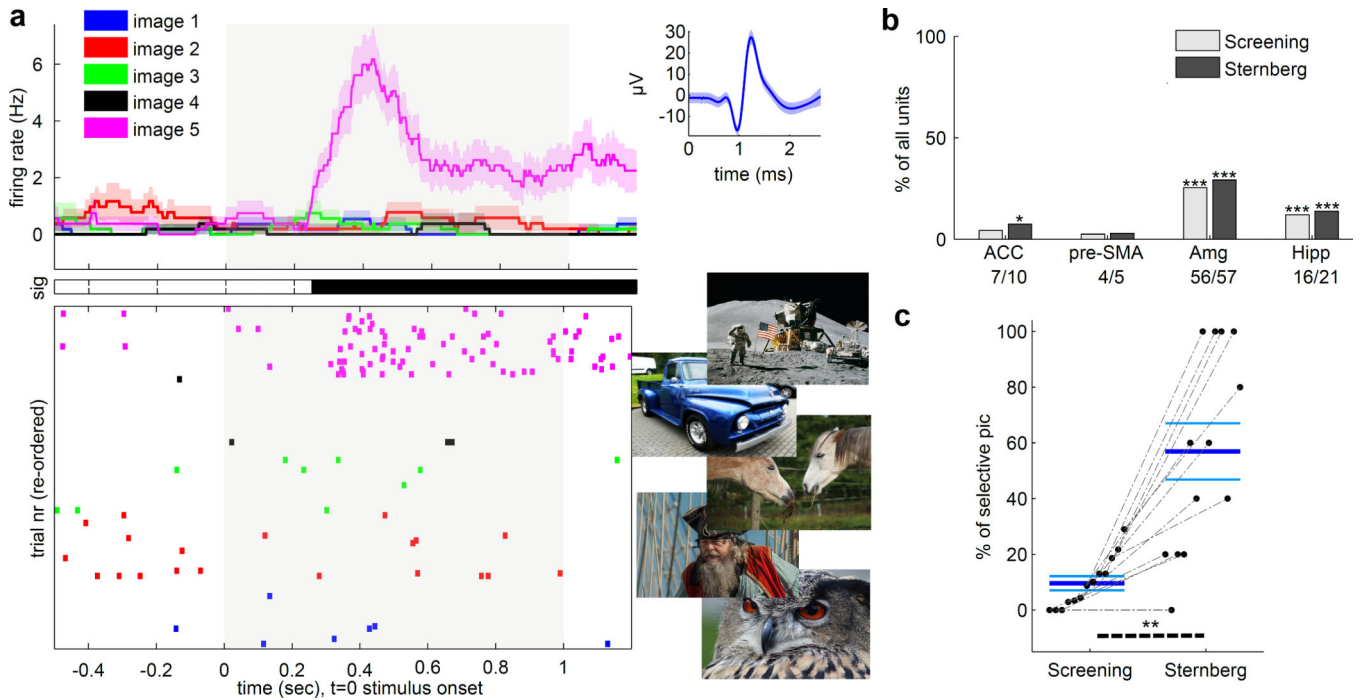


Figure 2. Stimulus-selective concept cells

(a) Example concept cell recorded from the amygdala. Upper panel shows the Post-stimulus time histogram (PSTH; binsize 200ms, stepsize 2ms). Colors denote the different images (shown on the right). Shaded areas represent \pm s.e.m across trials. Middle panel marks periods of significance (1×5 ANOVA; corrected for multiple comparisons using a cluster-size correction, see methods). Bottom panel shows raster with trials re-ordered according to image identity for plotting purposes only. Image onset is at $t=0$ (gray bar). The right inset shows the mean extracellular waveform \pm s.e.m of all spikes associated with this cell. (b) Percent of all recorded cells who qualified as concept cells in each area during the screening and WM (“Sternberg”) tasks. The numbers below each area label denotes the number of neurons associated with each bar. We test if the observed percentage is higher than that expected by chance by comparing with a null distribution estimated after scrambling the condition labels randomly (repeated 500 times; * denotes $p < 0.05$, ** $p < 0.01$, and *** denotes $p < 0.002$). (c) Percentage of images shown for which we observed at least one concept cell in both tasks ($p = 0.0014$). Each dashed line connects an individual session. This shows that the screening task successfully identified responsive neurons. Thick and thin blue lines represent the mean and \pm s.e.m, respectively.

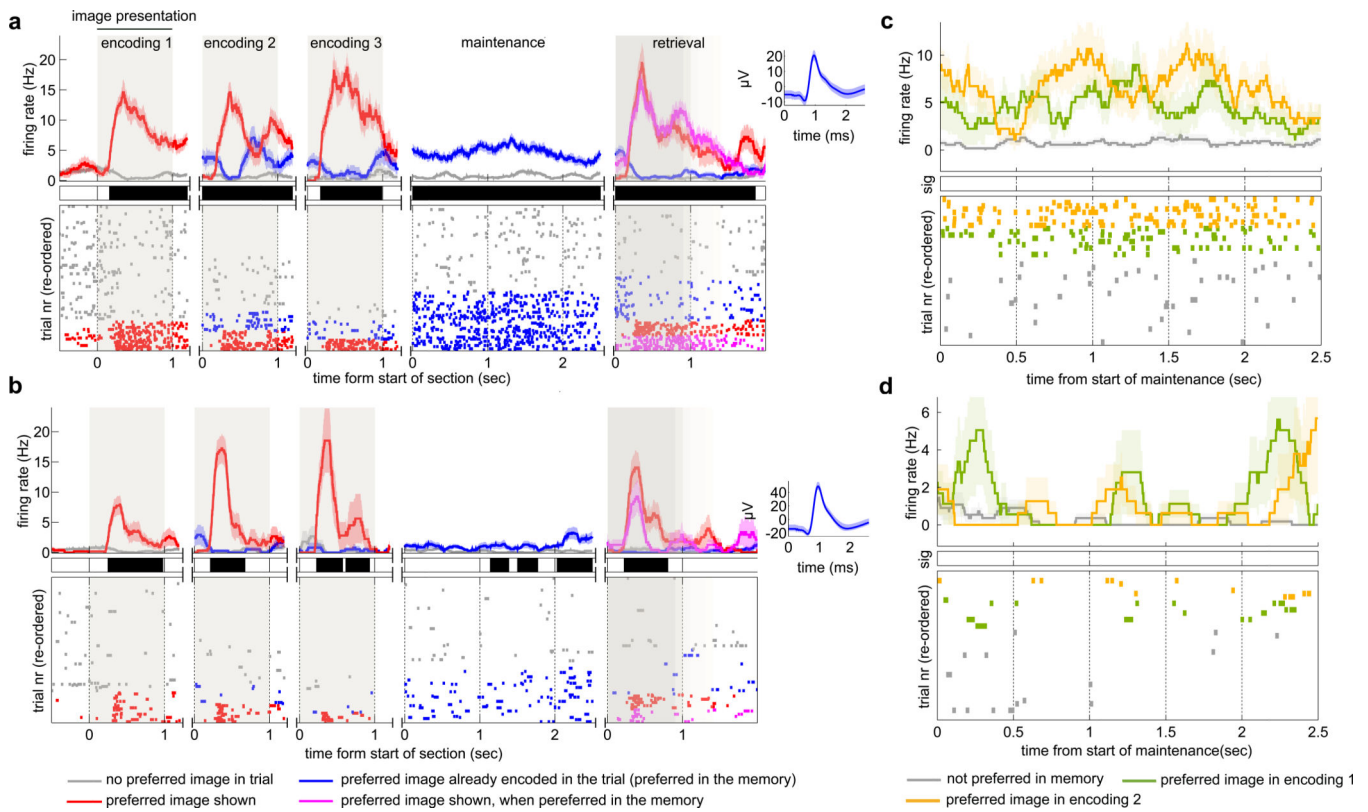


Figure 3. Concept cells are persistently active during WM maintenance

(a,b) Two example concept cells recorded from the amygdala (a) and hippocampus (b). For each, the upper panel shows the PSTH (binsize 200ms, stepsize 2ms). Shaded areas represent \pm s.e.m across trials. Middle panel marks periods of significance between preferred vs. not preferred stimuli (corrected for multiple comparisons using a cluster-size correction, see methods). Bottom panel shows raster with trials re-ordered according to condition for plotting purposes only. Both neurons show both visually evoked selective activity (red) and sustained activity (blue) during maintenance. Note how during maintenance, concept cells have elevated activity only when their preferred stimulus was held in memory (blue vs. gray). Also, note how the sustained activity (blue) was suppressed during encoding of the non-preferred image (i.e. encoding 3) when the preferred stimulus was already held in memory. (c,d) Maintenance activity of the same neurons shown in (a,b), but only for the subset of trials with load 2 (two items held in memory). There was no significant difference in activity between trials where the preferred image was shown first vs. second (encoding 1 or 2; middle panel marks periods of significance). See Fig. 4d for a similar analysis at the population level. See Fig. S3 for further single-cell examples.

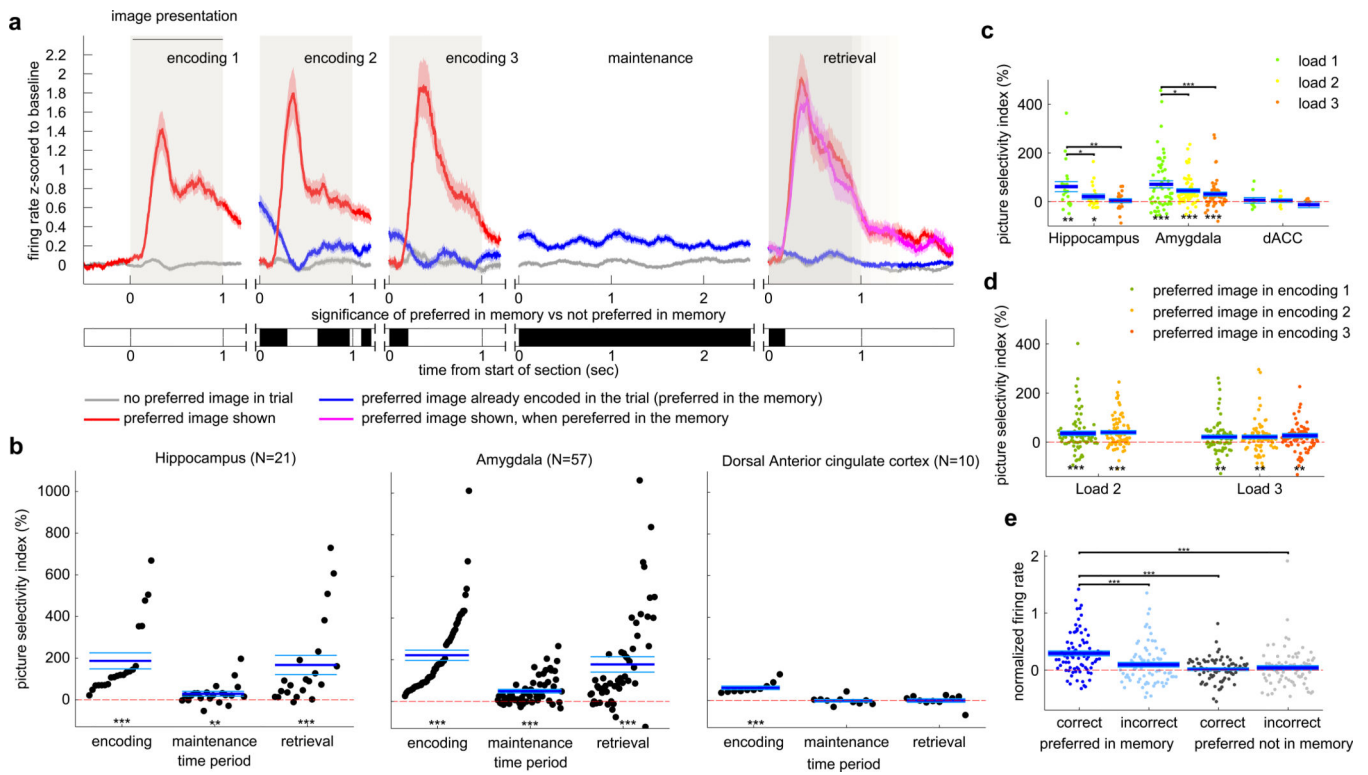


Figure 4. Population analysis of MTL concept cells

(a) Average firing rate of all concept cells identified in the amygdala ($n=57$) in the different phases of the task. Shaded areas represent \pm s.e.m across neurons. Gray vertical bars mark periods of time during which an image was on the screen. Bottom panel marks points of time during which the activity of the cells was significantly different between trials when a preferred image was in memory vs. when it was not (corrected for multiple comparisons based on cluster size, see methods). Colors mark different trials as indicated. For subplots a-d, only correct trials were used. (b) Picture selectivity index (PSI) during encoding, maintenance, and retrieval for all identified concept cells (each data point is one neuron; data points are sorted according to the encoding phase of the task). Neurons in both amygdala and hippocampus, but not dACC, maintained their selectivity throughout the task and showed persistent activity. Significance was computed against chance (PSI=0). (c) PSI for different load conditions indicates that neurons maintained persistent activity for loads 1–3 in amygdala and 1–2 in hippocampus, but not for dACC. (d) PSI for loads 2 and 3 as a function of whether the preferred image was shown first, second, or third during encoding. This shows that images which were shown directly before the maintenance period did not have greater selectivity (load 2 $P=0.702$; load 3 $P=0.873$). (e) Relationship between firing rate of concept cells in the MTL and behavior. The firing rate was significantly higher for correct compared to incorrect trials only when the preferred stimulus was held in memory. For (c–e), PSI and firing rate was calculated for the entire maintenance period. Throughout, * denotes $p<0.05$, ** $p<0.01$ and *** $p<0.001$ as estimated with permutation tests. Pre-SMA is not shown in this figure because we did not identify any concept cells in this area.

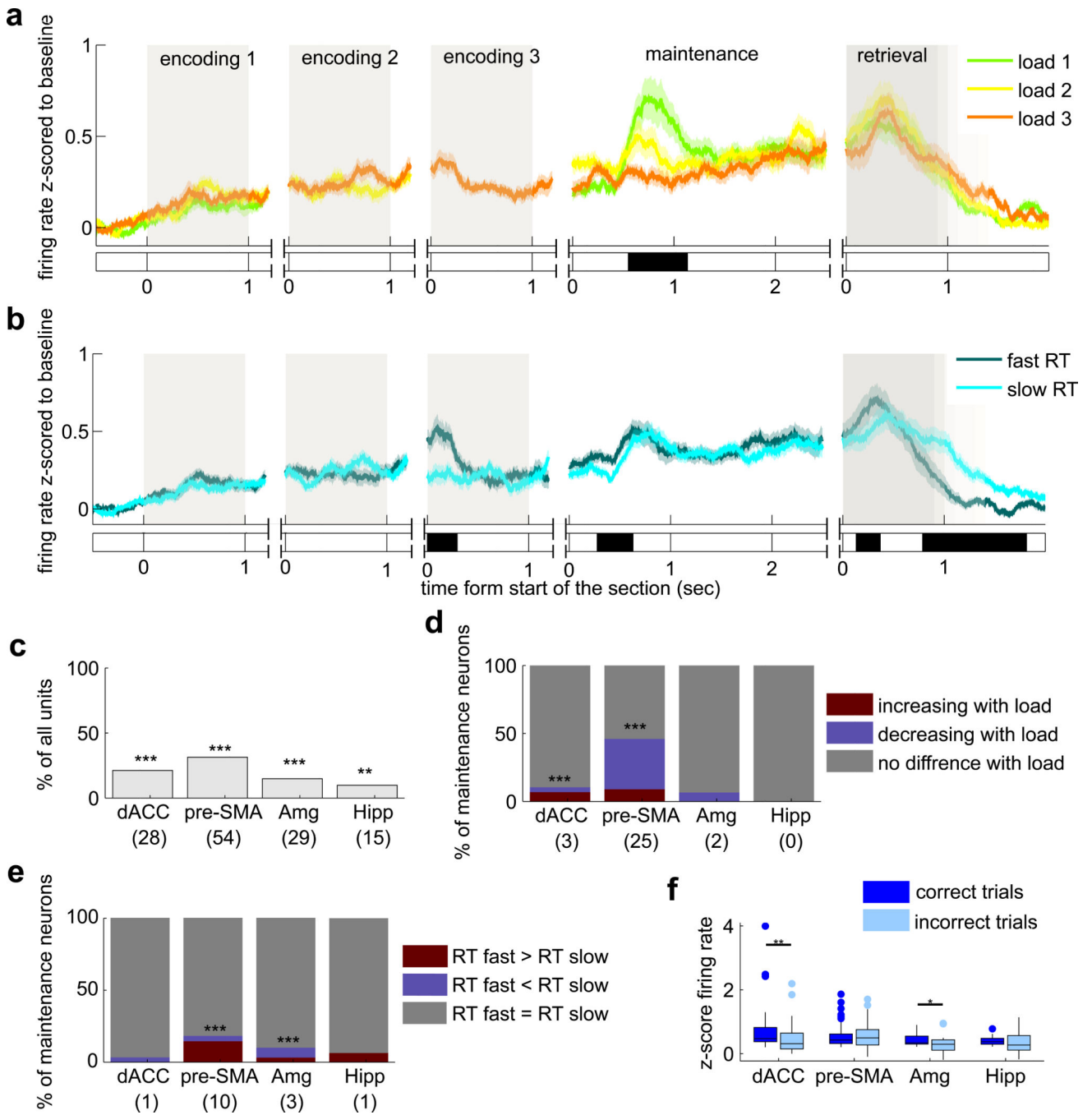


Figure 5. Persistent activity of maintenance cells in MFC

(a) Average firing rate of all maintenance neurons ($n=54$) in pre-SMA for different load conditions and (b) for trials that ended with a fast or slow response (median split of RT). Shaded areas represent \pm s.e.m across neurons. Gray vertical bars mark image presentation. Black bars indicate significance at $P<0.05$ of a 1×3 (top, permutation ANOVA) and 1×2 (bottom, permutation t-test) with load as dependent variable for (a) and response time for (b). Multiple comparisons were corrected for using a cluster-size approach (see methods). (c) Percentage of all recorded cells identified as maintenance neurons in each area. Numbers

below the area label denotes number of cells. The medial frontal areas (dACC, pre-SMA) contained significantly higher proportions of maintenance neurons ($\chi^2[1]=21.1$; $P=4.353e-6$) compared to areas in the MTL. (d) Percentage of maintenance neurons whose firing rate during maintenance differed as a function of load. Notably, in pre-SMA, 37% of cells decreased their firing rate as a function of load. (e) Percentage of maintenance neurons whose activity during maintenance differed as a function of response time. (f) The firing rate of maintenance neurons in dACC and amygdala differed as a function of whether stimuli were later remembered or forgotten. (c–e) Significance was assessed by comparing with a null distribution estimated using a bootstrap (see methods). For (f) we used permutation test. (f) Boxplot represents quartiles (25%, 75%), line is median, whiskers show range up to 1.5 times the interquartile range, and dots above whiskers show outliers. Throughout, * denotes $p < 0.05$, ** $p < 0.01$ and *** $p \leq 0.002$.

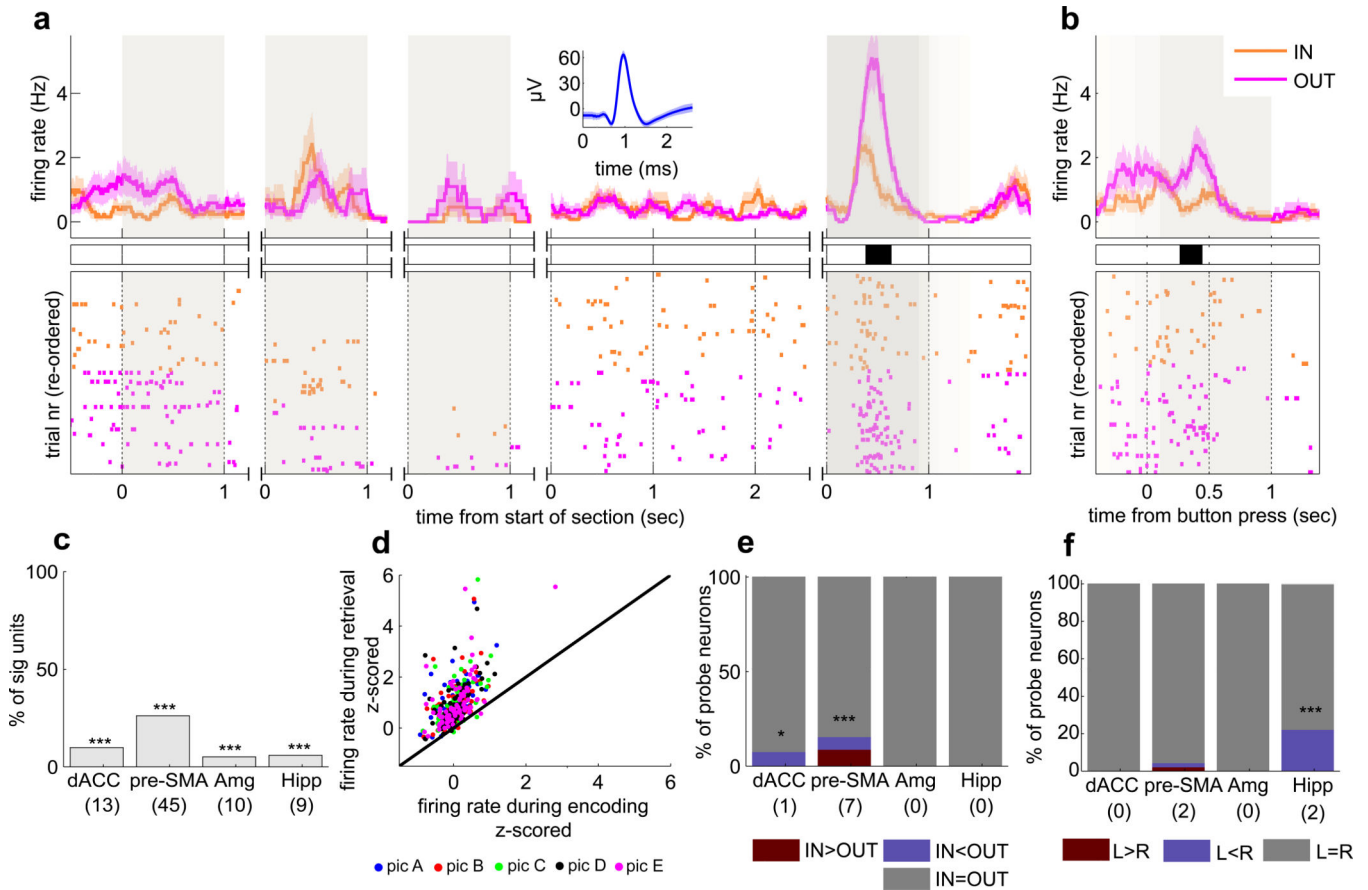


Figure 6. Probe neurons reflect WM-retrieval related evoked activity in MFC

(a) Firing rate of an example probe neuron recorded from the pre-SMA, shown separately for trials where the probe was held in (IN, cyan) or not (OUT, magenta) in memory. Upper panel shows the PSTH (binsize 200 bins, stepsize 2 ms, shaded areas represent \pm s.e.m). Middle panel marks points of time with a significant difference between IN and OUT trials, corrected for multiple comparisons using a cluster-size approach. Bottom panel shows raster with re-ordered trials. Gray vertical bars mark image presentation. (b) Same neuron as in (a), but aligned the response (button press was at 1.2 sec after image presentation) to button press. Note the much reduced peak response (1.44 Hz vs 0.45 Hz, permuted t-test: $P=0.005$). (c) Percentage of probe neurons in each area. Probe neurons were most prominent in pre-SMA, followed by dACC. (d) Probe neurons elevated their firing rate only during retrieval, but not encoding ($P=0.0002$, permuted t-test). (e–f) Percentage of probe neurons in each area whose firing rate during probe (–800–0 ms relative to button press) differed as a function of IN vs. OUT (e) or as a function of button press (f). Most cells showed no difference, i.e. they responded equally strongly (but selectively) to the probe stimulus. (c,e,f) Significance was assessed based on a null distribution estimated based on permuted labels. * denotes $p < 0.05$, ** $p < 0.01$ and *** $p < 0.001$.

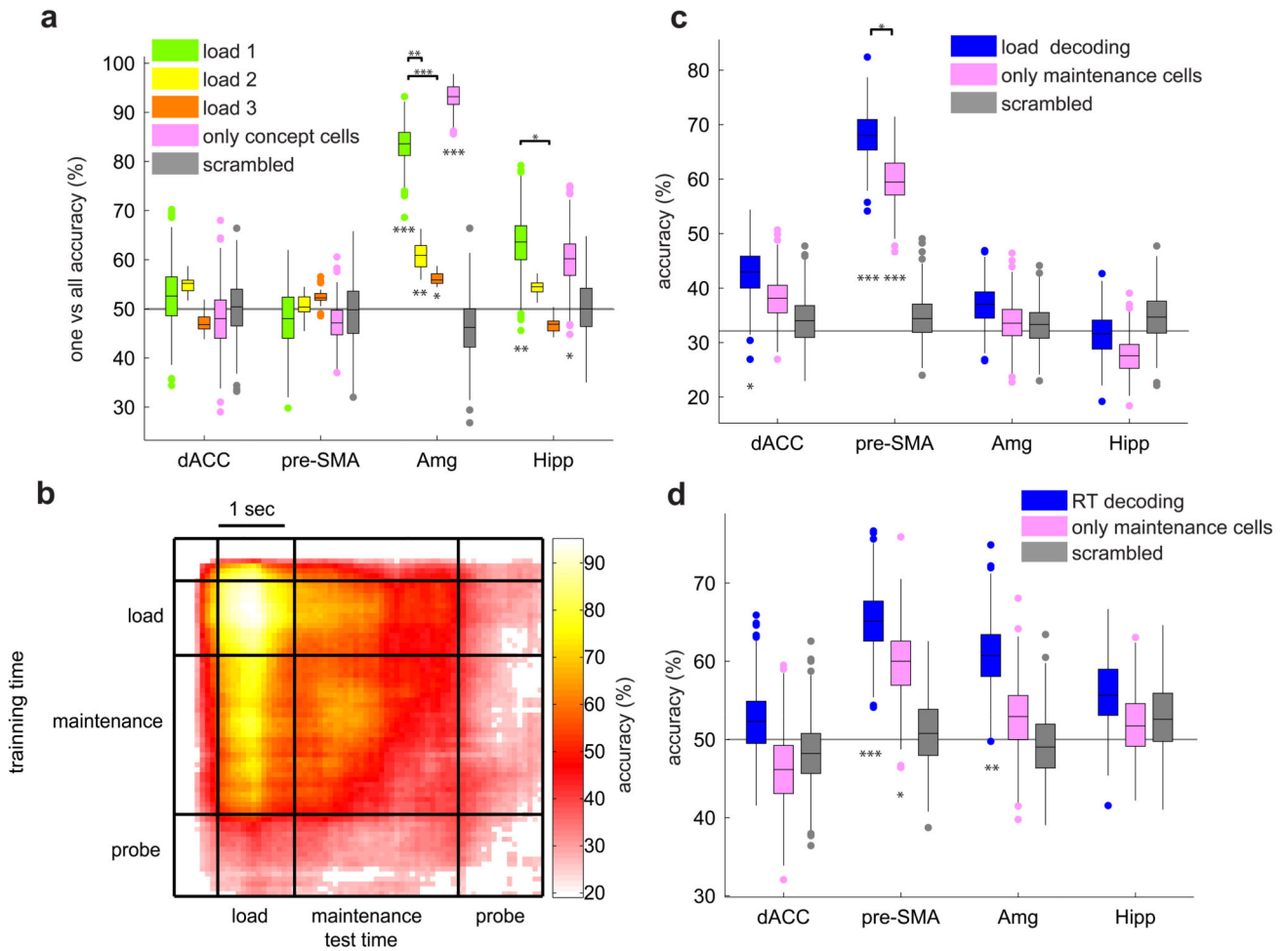


Figure 7. Population decoding from all recorded neurons during the maintenance period

(a) During maintenance, picture identity could be decoded from the activity of amygdala and hippocampus, but not dACC and pre-SMA, neurons. Decoding performance was maintained when only cells identified as concept cells were considered (magenta). One vs. all denotes average accuracy of decoders trained to distinguish between a given image and all the others (50% chance level). (b) Picture identity decoding using different time windows for training and testing. Shown is the test – retest decoding performance for load 1 trials (chance level is 20%). (c) The activity of neurons in the pre-SMA during maintenance was predictive of how many items were held in memory (load, 1–3). Decoding only from maintenance neurons (magenta) was sufficient. (d) The activity of neurons in the pre-SMA and amygdala during maintenance was predictive of later response speed. * denotes significance at $p < 0.05$, ** $p < 0.01$ and *** $p \leq 0.002$. Markers below bars indicate significance vs. chance performance, estimated by randomly scrambled labels. Significance of pairwise tests was estimated by comparisons with a null distribution of the same differences estimated from decoders trained on data with randomly scrambled labels. (a,c,d) Boxplots represent quartiles (25%, 75%), line indicates the median, whiskers show range up to 1.5 times the interquartile range, and dots above whiskers show outliers.

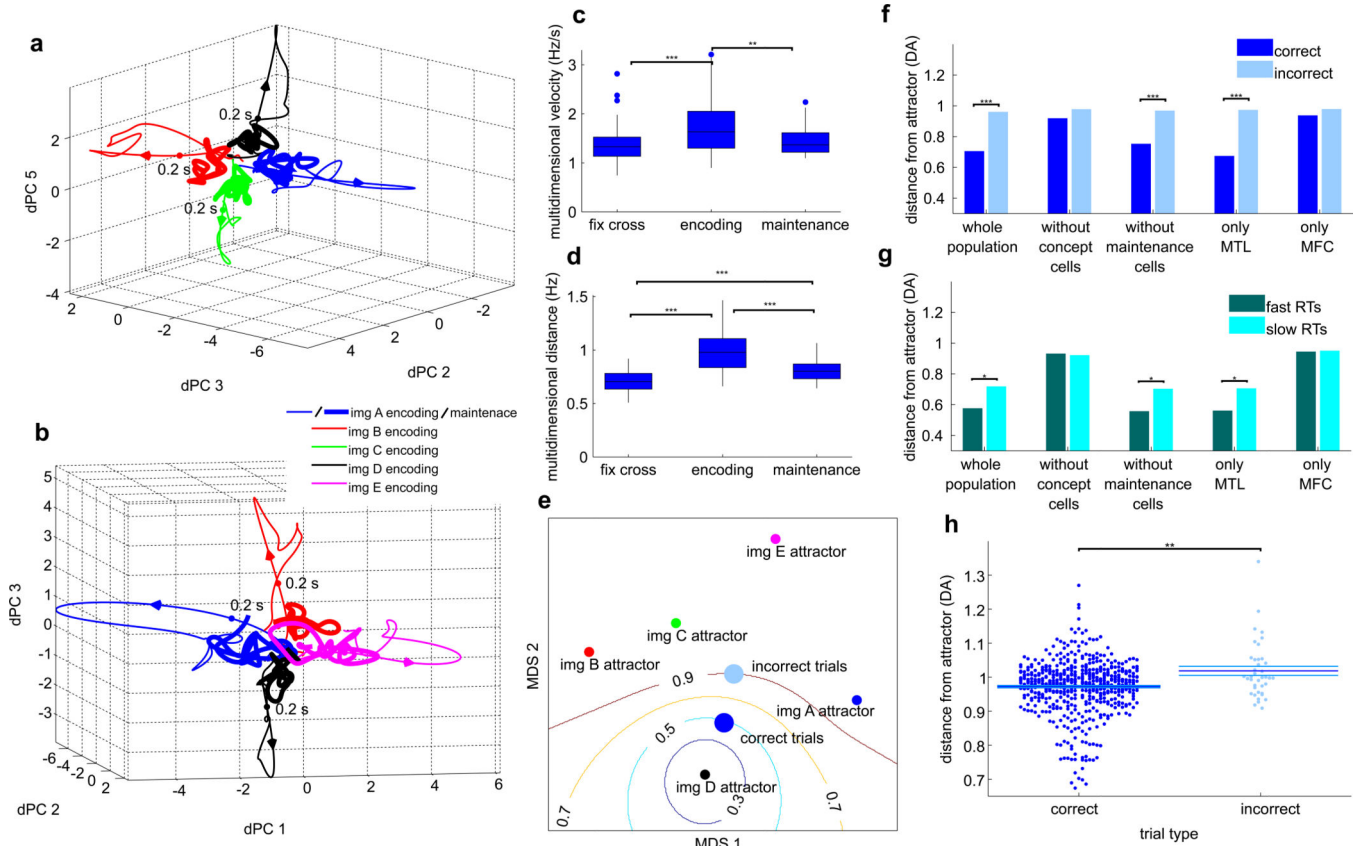


Figure 8. Persistent activity during maintenance forms attractors

(a) Illustration of the mean trajectories in neuronal state space formed by the three demixed principal components (dPCs) associated with picture identity during encoding (thin line) and maintenance (thick line). The dot indicates the point of time 200 ms after image onset and the arrow indicates the direction of change. Colors mark different images (only 4 of the total 5 are shown for clarity). (b) Different view of (a). (c) Multidimensional velocity of the population in the different phases of the task (shown for all load 1 trials). The velocity during maintenance was significantly slower compared to encoding and was not significantly different from that during baseline. (d) Multidimensional pairwise distance between all possible pairs of attractors during maintenance (load 1). The distance during maintenance was significantly larger compared to that during baseline ($P=0.0005$). Together (c,d) are indicative of attractors. The significance of the population-metrics shown in (c–d) was computed by randomly subsampling a subset of trials and neurons (see methods). (e) Schematic representation of the distances between the attractors (attractors are defined based on correct load 1 trials) for each image (small filled dots) and the average position in state space for image D for two behaviors: remembered (correct) and forgotten (incorrect) computed for all loads separately and averaged. This representation was determined based on multidimensional scaling of the state space (see methods). Isolines depict areas of equal distance from the attractor for image D. (f) The distance to the attractor (DA) was significantly smaller for correct compared to incorrect trials only when concept cells were part of the population. Note that $DA < 1$ indicates that the trajectory is closer to the correct attractor than all the other attractors. (g) The distance to the attractor (DA) corresponding to

the remembered image was indicative of the speed of the response. This relationship was observed only for concept cells. (h) Distance from attractor predicts performance on individual trials. Each dot is one trial. (c,d) Boxplots represent quartiles (25%, 75%), line indicates the median, whiskers show range up to 1.5 times the interquartile range, and dots above whiskers show outliers. * denotes $p < 0.05$, ** $p < 0.01$ and *** denotes $p \leq 0.002$.