

# 1 Perception and memory have distinct spatial tuning 2 properties in human visual cortex

3 Serra E. Favila<sup>1\*+</sup>, Brice A. Kuhl<sup>2</sup>, and Jonathan Winawer<sup>1</sup>

4 <sup>1</sup>Department of Psychology, New York University, New York, NY, 10003

5 <sup>2</sup>Department of Psychology, University of Oregon, Eugene, Oregon, 97403

6 \*Contact: sef2177@columbia.edu

7 +Current address: Department of Psychology, Columbia University, New York, NY, 10027

## 8 Abstract

9 Reactivation of earlier perceptual activity is thought to underlie long-term memory recall. Despite evidence for  
10 this view, it is unknown whether mnemonic activity exhibits the same tuning properties as feedforward perceptual  
11 activity. Here, we leveraged population receptive field models to parameterize fMRI activity in human visual cortex  
12 during spatial memory retrieval. Though retinotopic organization was present during both perception and memory,  
13 large systematic differences in tuning were also evident. Notably, whereas there was a three-fold decline in spatial  
14 precision from early to late visual areas during perception, this property was entirely abolished during memory  
15 retrieval. This difference could not be explained by reduced signal-to-noise or poor performance on memory  
16 trials. Instead, by simulating top-down activity in a network model of cortex, we demonstrate that this property is  
17 well-explained by the hierarchical structure of the visual system. Our results provide insight into the computational  
18 constraints governing memory reactivation in sensory cortex.

20 **Keywords:** episodic memory, spatial memory, reinstatement, reactivation, visual cortex, population receptive field,  
21 hierarchical model

## 22 Introduction

23 Episodic memory retrieval allows humans to bring to mind the details of a previous experience. This process is  
24 hypothesized to involve reactivating sensory activity that was evoked during the initial event (James, 1890; Hebb,  
25 1968; Damasio, 1989; McClelland et al., 1995). For example, remembering a friend's face is thought to involve  
26 reactivating neural activity that was present when seeing that face. There is considerable evidence from human  
27 neuroimaging demonstrating that the same visual cortical areas active during perception are also active during  
28 imagery and long-term memory retrieval (Kosslyn et al., 1995; O'Craven & Kanwisher, 2000; Wheeler et al., 2000;  
29 Slotnick et al., 2005; Polyn et al., 2005; Kuhl et al., 2011; Bosch et al., 2014; Waldhauser et al., 2016; Lee et al.,  
30 2018; Bone et al., 2018). These studies have found that mnemonic activity in early visual areas like V1 reflects  
31 the low-level visual features of remembered stimuli, such as spatial location and orientation (Kosslyn et al., 1995;  
32 Thirion et al., 2006; Bosch et al., 2014; Naselaris et al., 2015; Sutterer et al., 2019). Likewise, category-selective  
33 activity in high-level visual areas like FFA and PPA is observed when subjects remember or imagine faces and  
34 houses (O'Craven & Kanwisher, 2000; Polyn et al., 2005). The strength and pattern of visual cortex activity has  
35 been associated with retrieval success in memory tasks (Kuhl et al., 2011, 2013; Gordon et al., 2014), suggesting  
36 that cortical reactivation is relevant for behavior.

37 These studies, and many others, have established similarities between the neural substrates of visual perception  
38 and visual memory. However, relatively less attention has been paid to identifying and explaining *differences*  
39 between activity patterns evoked during perception and memory. In the present work, we asked the following  
40 question: which properties of stimulus-driven activity are reproduced in visual cortex during memory retrieval and  
41 which are not? The extreme possibility—that all neurons in the visual system produce identical responses when  
42 perceiving vs remembering a given stimulus—can likely be rejected. Early studies demonstrated that sensory

43 responses were reduced during memory retrieval relative to perception (Wheeler et al., 2000), and perception and  
44 memory give rise to distinct subjective experiences. A more plausible proposal is that visual memory functions as a  
45 "weak" version of feedforward perception (Pearson et al., 2015; Pearson, 2019), with memory activity organized in  
46 the same fundamental way as perceptual activity, but with reduced signal-to-noise. This hypothesis is consistent  
47 with informal comparisons between perception and memory BOLD amplitudes and data suggesting that visual  
48 imagery produces similar behavioral effects to weak physical stimuli in many tasks (Ishai & Sagi, 1995; Pearson  
49 et al., 2008; Tartaglia et al., 2009; Winawer et al., 2010). A third possibility is that memory reactivation differs from  
50 stimulus-driven activation in predictable and systematic ways beyond signal-to-noise. Such differences could arise  
51 due to a change in the neural populations recruited, a change in those populations' response properties, or a  
52 systematic loss of information during sensory encoding or post-sensory processing.

53 One way to adjudicate between these possibilities is to make use of models from visual neuroscience that  
54 quantitatively parameterize the relationship between stimulus properties and the BOLD response. In the spatial  
55 domain, population receptive field models (pRF) define a 2D receptive field that transforms stimulus position on the  
56 retina to a voxel's BOLD response (Dumoulin & Wandell, 2008; Wandell & Winawer, 2015). These models are based  
57 on well-understood physiological properties of the primate visual system and account for a large amount of variance  
58 in the BOLD signal observed in human visual cortex during perception (Kay et al., 2013b). Using these models to  
59 quantify memory-evoked activity in the visual system offers the opportunity to precisely model the properties of  
60 memory reactivation in visual cortex and their relationship to visual activation. In particular, the fact that pRF models  
61 describe neural activity in terms of stimulus properties may aid in interpreting differences between perception and  
62 memory activity patterns by projecting these differences onto a small number of interpretable physical dimensions.

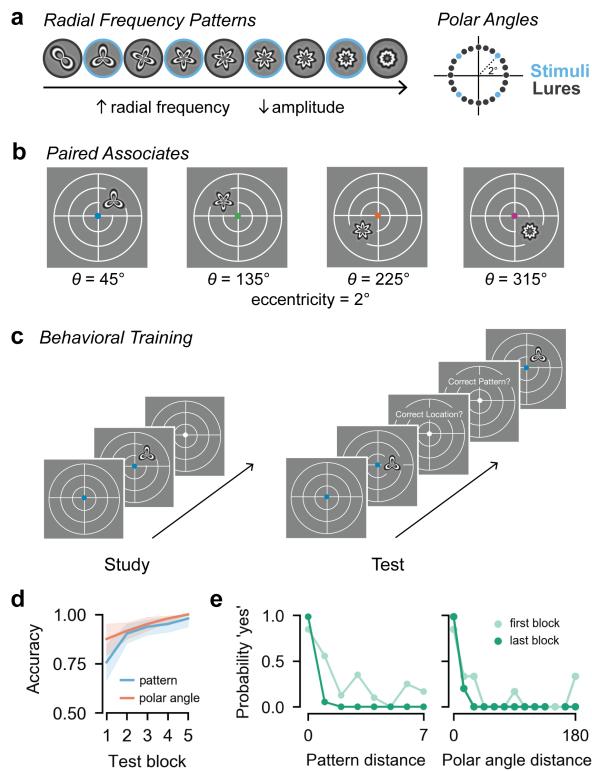
63 Here, we used pRF models to characterize the spatial tuning properties of mnemonic activity in human visual  
64 cortex. We first trained human subjects to associate spatially localized stimuli with colored fixation cues. We then  
65 measured stimulus-evoked and memory-evoked activity in visual cortex using fMRI. Separately, we fit pRF models  
66 to independent fMRI data, which allowed us to estimate receptive field location and size within multiple visual field  
67 maps for each subject. Using pRF-based analyses, we quantified the location, amplitude, and precision of neural  
68 activity within these visual field maps during perception and memory retrieval. Finally, we explored the cortical  
69 computations that could account for our observations by simulating neural responses using a stimulus-referred pRF  
70 model and a hierarchical model of neocortex.

## 71 Results

### 72 Behavior

73 Prior to being scanned, subjects participated in a behavioral training session. During this session, subjects  
74 learned to associate four colored fixation dot cues with four stimuli. The four stimuli were unique radial frequency  
75 patterns presented at 45, 135, 225, or 315 degrees of polar angle and 2 degrees of eccentricity (Fig. 1a,b).  
76 Subjects alternated between study and test blocks (Fig. 1c). During study blocks, subjects were presented with the  
77 associations. During test blocks, subjects were presented with the cues and had to detect the associated stimulus  
78 pattern and polar angle location among similar lures (Fig. 1a,c; see Methods). All subjects completed a minimum of  
79 4 test blocks (mean = 4.33, range = 4–5), and continued the task until they reached 95% accuracy. Subjects' overall  
80 performance improved over the course of training session (Fig. 1d). In particular, subjects showed improvements in  
81 the ability to reject similar lures from the first to the last test block (Fig. 1e).

82 After subjects completed the behavioral training session, we collected fMRI data while subjects viewed and  
83 recalled the stimuli (Fig. 2a). During fMRI perception runs, subjects fixated on the central fixation dot cues and  
84 viewed the four stimuli in their learned spatial locations. Subjects performed a one-back task to encourage covert  
85 attention to the stimuli. Subjects were highly accurate at detecting repeated stimuli (mean = 86.9%, range =  
86 79.4%–93.2%). During fMRI memory runs, subjects fixated on the central fixation dot cues and recalled the  
87 associated stimuli in their spatial locations. On each trial, subjects made a judgment about the subjective vividness  
88 of their memory. Subjects reported that they experienced vivid memory on an average of 89.8% of trials (range:  
89 72.4%–99.5%), weak memory on 8.9% of trials (0.5%–25.0%), and no memory on 1.2% of trials (0.5%–2.6%).

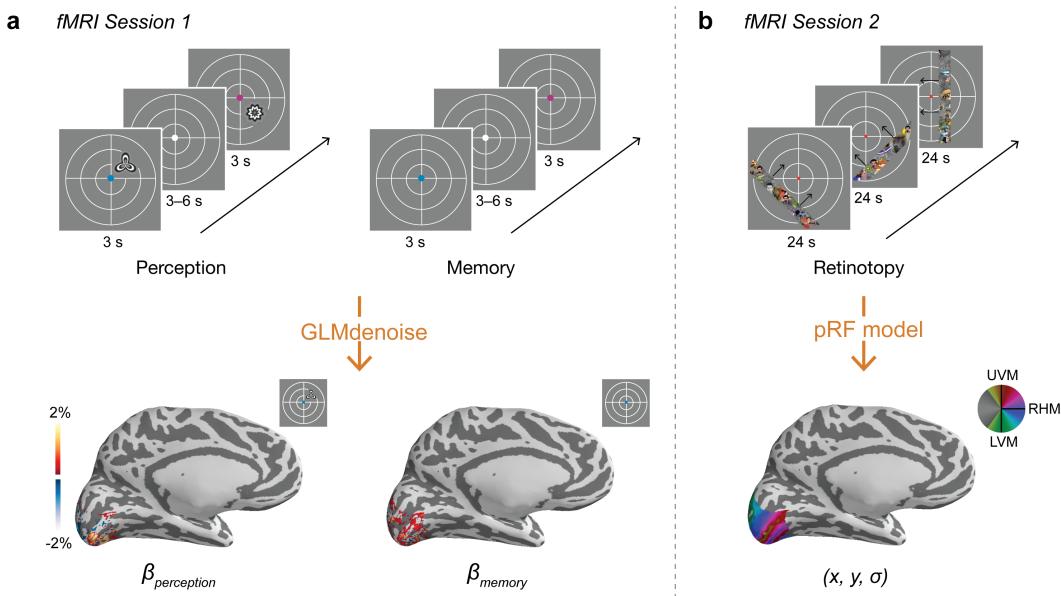


**Figure 1. Stimuli and behavioral training.** (a) The four radial frequency patterns and polar angle locations used in the fMRI experiment are outlined in blue. The intervening patterns and locations were used as lures during the behavioral training session. (b) Immediately prior to the scan, subjects learned that each of four colored fixation dot cues was associated with a unique radial frequency pattern that appeared at a unique location in the visual field. (c) During training, subjects alternated between study and test blocks. During study blocks, subjects were presented with the associations while maintaining central fixation. During test blocks, subjects were presented with the cues followed by test probes while maintaining central fixation. Subjects gave yes/no responses to whether the test probe was presented at the target polar angle and whether it was the target pattern. (d) Accuracy of pattern and polar angle responses improved over the course of the training session. Lines indicate average accuracy across subjects. Shaded region indicates 95% confidence interval. (e) Memory performance became more precise from the first to the last test block. During the first block, false alarms were high for stimuli similar to the target. These instances decreased by the last test block. Dots indicate probability of a 'yes' response for all trials and subjects in either the first or last block. The x axis is organized such that zero corresponds to the target and increasing values correspond to lures more dissimilar to the target.

## Memory reactivation is spatially organized

We used a GLM to estimate the BOLD response evoked by seeing and remembering each of the four spatially localized stimuli (Fig. 2a; see Methods). Separately, each subject participated in a retinotopic mapping session. We fit pRF models to these data to estimate pRF locations ( $x, y$ ) and sizes ( $\sigma$ ) in multiple visual areas (Fig 2b). To more easily compare perception- and memory-evoked activity across visual areas, we transformed these responses from cortical surface coordinates into visual field coordinates using the pRF parameters. For each subject, ROI, and stimulus, we plotted the amplitude of the evoked response at the visual field position ( $x, y$ ) estimated by the pRF model (Fig. 3a). We then interpolated these values over 2D space, z-scored the values, rotated all stimulus responses to the same polar angle, and averaged across stimuli and subjects (see Methods). These plots are useful for comparison across regions because they show the organization of the BOLD response in a common space that is undistorted by the size and magnification differences present in cortex.

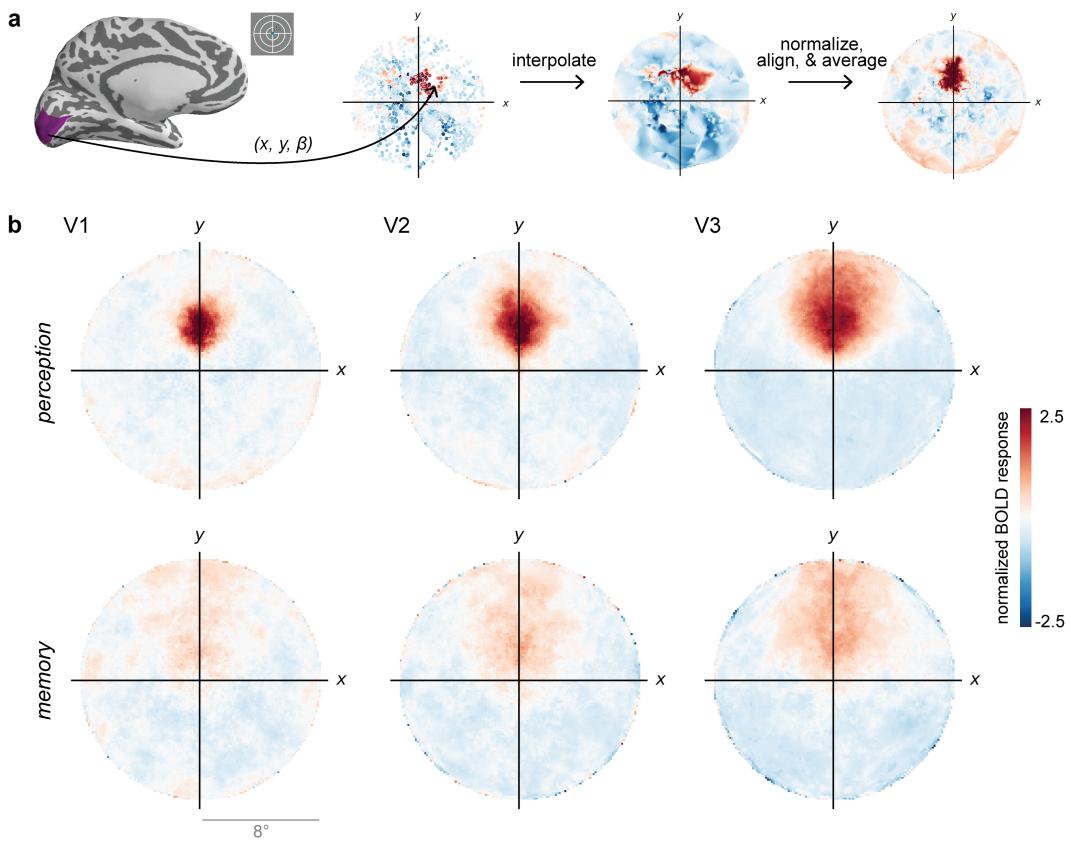
We generated these visual field plots for V1, V2, and V3 as an initial way to visualize the evoked responses during perception and memory. Readily apparent is the fact that stimulus-evoked responses during perception were robust and spatially-specific (Fig. 3b, top). The spatial spread of perceptual responses increased from V1 to V3, consistent with estimates of increasing receptive field size in these regions (Wandell & Winawer, 2015; Kay et al., 2013b). While the memory responses were weaker and more diffuse, they were also spatially organized, with peak activity in the same location as the perception responses (Fig. 3b, bottom).



**Figure 2. fMRI task design and measurements.** (a) Following training, subjects participated in two tasks while being scanned. During perception runs, subjects viewed the colored fixation dot cues and associated stimuli while maintaining central fixation. Subjects performed a one-back task on the stimuli to encourage covert attention to each stimulus. During memory runs, subjects viewed only the cues and recalled the associated stimuli while maintaining central fixation. Subjects made a judgment about the vividness of their memory (vivid, weak, no memory) on each trial. We used the perception and memory fMRI time series to perform a GLM analysis that estimated the response evoked by perceiving and remembering each stimulus for each vertex on the cortical surface. Responses in visual cortex for an example subject and stimulus are shown at bottom. (b) In a separate fMRI session on a different day, subjects participated in a retinotopic mapping session. During retinotopy runs, subjects viewed bar apertures embedded with faces, scenes, and objects drifting across the visual field while they maintained central fixation. Subjects performed a color change detection task on the fixation dot. We used the retinotopy fMRI time series to solve a pRF model that estimated the receptive field parameters for each vertex on the cortical surface. A polar angle map is plotted for an example subject at bottom.

107 We quantified these initial observations. Because our stimulus locations were isoeccentric, we reduced our  
 108 responses to variance along one spatial dimension: polar angle. To do this, we restricted our ROIs to surface  
 109 vertices with pRF locations near the stimulus eccentricity, rotated stimuli to a common polar angle, normalized the  
 110 responses, and averaged across stimuli and subjects (see Methods). We then plotted the group average BOLD  
 111 response in bins of polar angle distance from the stimulus (Fig. 4a). We generated these polar angle response  
 112 functions for V1–V3 and for three mid-level visual areas: hV4, LO, and V3ab (Fig. 4b). To capture the pattern of  
 113 positive and negative BOLD responses we observed, we fit the average data in each ROI with a difference of two  
 114 von Mises distributions, where both the positive and the negative von Mises were centered at the same location.  
 115 Visualizing the data and the von Mises fits (Fig. 4b), it's clear that both perception and memory fits show a peak at  
 116 0°, or the true location of the stimulus, in every region.

117 To formally test this, we calculated bootstrapped confidence intervals for the location parameter of the von Mises  
 118 distributions by resampling subjects with replacement (see Methods). We then compared the accuracy and reliability  
 119 of location parameters between perception and memory (Fig. 4c, left). As expected, location parameters derived  
 120 from perception data were highly accurate. 95% confidence intervals for perception location parameters overlapped  
 121 0° of polar angle, or the true stimulus location, in all ROIs. These confidence intervals spanned only 7.0° on average  
 122 (range: 3.9°–9.5°), demonstrating that there was low variability in location accuracy across subjects in every ROI.  
 123 Critically, memory parameters were also highly accurate, with confidence intervals overlapping 0° in every ROI (Fig.  
 124 4c, left). Thus, in every visual area measured, the spatial locations of the remembered stimuli could be accurately  
 125 estimated from mnemonic activity. Memory confidence intervals spanned 17.6° on average (range = 11.0°–21.3°),  
 126 indicating that location estimates were somewhat less reliable during memory during perception. However, even  
 127 the widest memory confidence interval spanned only 21.3°. This is far less than the 90° separating each stimulus  
 128 location, suggesting that there was no confusability between stimuli present in distributed memory activity. Because  
 129 both perception and memory location parameters were highly accurate, and because differences in reliability were  
 130 relatively small, there was no overall difference between perception and memory in the estimated location of peak

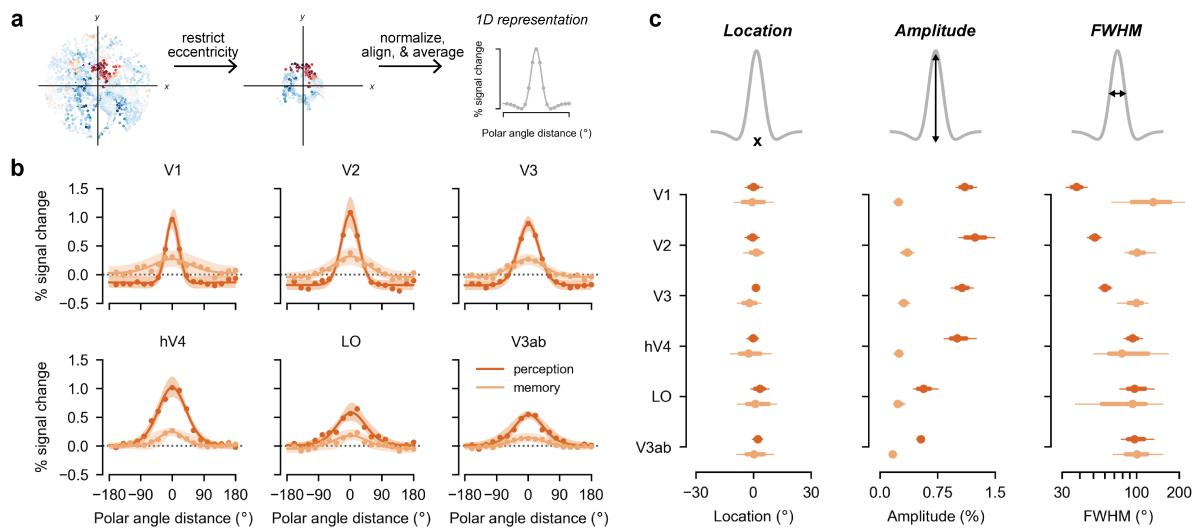


**Figure 3. Perception and memory activity in visual field coordinates.** (a) For a given subject, ROI, and stimulus, we plotted the perception- or memory-evoked response ( $\beta$ ) in the visual field position estimated by the pRF model ( $x, y$ ). We then interpolated over 2D space and z-scored the responses. We rotated these representations by the polar angle location of the stimulus so that they aligned on the upper vertical meridian, and then averaged over stimuli. This procedure produces an average activation map in visual field coordinates for each ROI and subject. This map is plotted for V1 in an example subject, at right. (b) Plots of perception-evoked and memory-evoked activity, averaged across all subjects, from V1, V2, and V3. These plots reproduce known features of spatial processing during perception, such as increasing receptive field size from V1–V3. They also qualitatively demonstrate that perceptual activity is not perfectly reproduced during memory retrieval but that some retinotopic organization is preserved.

activity (main effect of perception/memory:  $\beta = 0.14$ , 95% CI = [-6.72, 6.14],  $p = 0.94$ ; Fig. 4c, left). These results provide strong evidence that memory-triggered activity in human visual cortex is spatially organized within known visual field maps, as it is during visual perception. These findings support prior reports of retinotopic activity during memory and imagery (Kosslyn et al., 1995; Slotnick et al., 2005; Thirion et al., 2006), but provide more quantitative estimates of this effect.

### Amplitude and precision differ between perceptual and mnemonic activity

Aspects of perception and memory responses other than the peak location differed considerably. First, memory responses were lower in amplitude than perception responses (Fig. 4b). To quantify this observation, we derived a measure of amplitude from the difference of von Mises functions fit to our data (see Methods). We also computed bootstrapped confidence intervals for this amplitude metric, following the prior analysis. We then compared these estimates between perception and memory. First, response amplitudes for perception data were higher than for the memory data (main effect of perception/memory:  $\beta = 0.95$ , 95% CI = [0.80, 1.13],  $p = 0.013$ ; Fig. 4c, middle). The average amplitude during perception was 0.92% signal change, and the average amplitude during memory was 0.26% signal change. Amplitude confidence intervals for perception and memory did not overlap in any ROI, indicating that these differences were highly significant in each region. Critically, the fact the perception amplitudes were larger than memory amplitudes does not imply that memory responses were at baseline. In fact, 95% confidence intervals for memory amplitudes did not overlap with zero in any region (Fig. 4c, middle),



**Figure 4. Perception and memory have shared and distinct activation features.** (a) We created 1D polar angle response functions by restricting data to eccentricities near the stimulus, aligning stimuli to a common polar angle, and averaging responses into polar angle distance bins. A difference of two von Mises distributions was fit to the group average response. Responses in cortical areas that have pRFs near the stimulus position are plotted at  $x = 0$ . (b) Polar angle response functions, averaged across all subjects and stimuli, are plotted separately for perception and memory. Dots represent average data across all stimuli and subjects. Lines represent the fit of the difference of two von Mises distributions to the average data, and shading represents the 95% confidence interval around this fit. While the peak location of the response is shared across perception and memory, there are clear differences in the amplitude and width of the responses. (c) Bootstrapped 68% (thick lines) and 95% confidence intervals (thin lines) for the location, amplitude and FWHM of the difference of von Mises fits are plotted to quantify the responses. In all ROIs, the peak location of the response is equivalent during perception and memory (at 0°, the stimulus location), while the amplitude of the response is reliably lower during memory than during perception. The FWHM of the response increases across ROIs during perception but not during memory, resulting in highly divergent FWHM for perception and memory in early visual areas.

demonstrating that responses were significantly above baseline in all areas measured. These results demonstrate that the amplitude of spatially-organized activity in visual cortex is attenuated (but present) during memory retrieval.

Second, memory responses were wider than perception responses (Fig. 4b). We operationalized the *precision* of perception and memory responses by computing the full width at half maximum (FWHM) of the difference of von Mises fit to our data and by generating confidence intervals for this measure. Note that FWHM is not sensitive to the overall scale of the response function: a perception response function rescaled to have the same amplitude as the memory response function will have an unchanged FWHM. On average, FWHM during perception was significantly smaller than during memory (main effect of perception/memory:  $\beta = -75.2$ , 95% CI = [-138.5, -33.1],  $p = 0.0002$ ; Fig. 4c, right). However, these differences were not equivalent across ROIs (perception/memory  $\times$  ROI interaction:  $\beta = 18.8$ , 95% CI = [5.78, 35.5],  $p = 0.021$ ). Specifically, perception FWHM increased moving up the visual hierarchy (main effect of ROI:  $\beta = 13.3$ , 95% CI = [10.3, 20.6],  $p = 0.0056$ ), indicating increased width or *decreased precision* in later visual areas compared to early visual areas (Fig. 4c, right). For example, V1 had the narrowest (most precise) response during perception, with an average FWHM of 38.0°(95% CI: [32.0°, 45.0°]), while V3ab had the widest responses during perception, with a FWHM of 97.0°(95% CI: [78.0°, 132.5°]). This increasing pattern follows previously described increases in population receptive field size in these regions (Wandell & Winawer, 2015; Kay et al., 2013b). Note that a separate question, not addressed here, is the precision with which the stimulus can be decoded from a representation, which is not necessarily related to receptive field size.

Strikingly, this pattern of increasing FWHM from early to late visual areas was abolished during memory (main effect of ROI:  $\beta = -5.49$ , 95% CI = [-18.7, 8.41],  $p = 0.20$ ; Fig. 4c, right). For areas V1–hV4, the regions we can sort hierarchically with the most confidence, the pattern across ROIs trended toward being reversed, with the widest responses observed in the *earliest* areas (main effect of ROI:  $\beta = -15.7$ , 95% CI = [-39.8, 12.4],  $p = 0.083$ ). These data demonstrate that fundamental aspects of spatial processing commonly observed during perception do not generalize to memory-evoked responses. Interestingly, the interaction between perception/memory and ROI yielded highly divergent perception and memory responses in the earliest visual areas but equivalent responses in the latest

172 visual areas (Fig. 4c, right). For example, V1 responses during memory had an average FWHM of 131.0° (95%  
173 CI: [66.9°, 225.0°]), and were thus 3.45 times wider than V1 responses during perception. In V2 and V3, memory  
174 FWHM exceeded perception FWHM by an average of 1.98 times and 1.67 times, respectively. Unlike in V1-V3,  
175 confidence intervals for perception and memory were highly overlapping in hV4, LO, and V3ab (Fig. 4c, right). In  
176 these later areas, memory responses were only 0.84–1.04 times wider during memory than during perception.  
177 These data raise the interesting possibility that later stages of perceptual processing serve as a bottleneck on  
178 mnemonic activity precision. Taken together, these results provide evidence for reliable and striking differences in  
179 the precision of perception and memory activity across different levels of the visual system. More broadly, these  
180 findings indicate that there are fundamentally different constraints on the properties of feedforward perceptual  
181 activity and top-down mnemonic activity in human visual cortex.

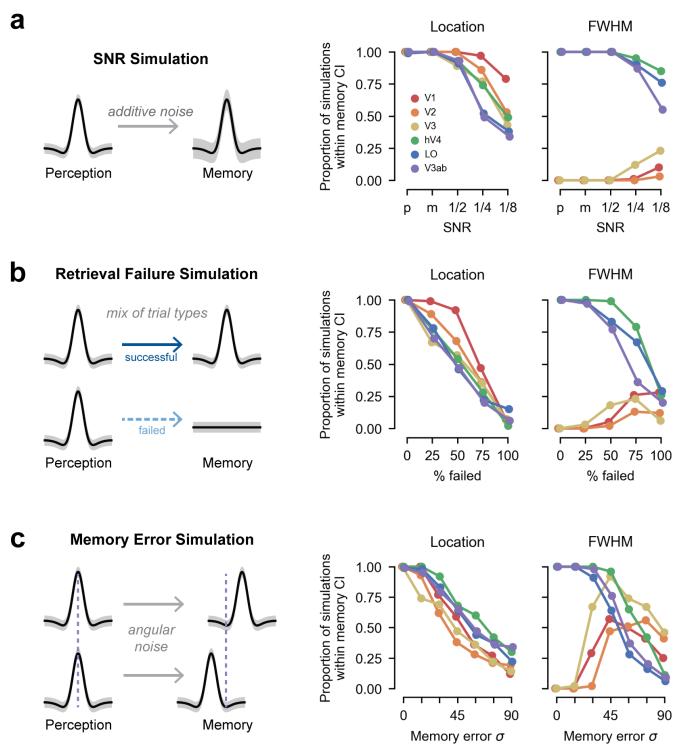
## 182 Differences between perception and memory responses are not explained by noise

183 One important consideration in interpreting our results is whether the differences we observed between perception  
184 and memory could be caused by differences in noise. For example, is it possible that perception and memory  
185 responses were actually equivalent other than noise level, but due to greater trial-to-trial noise, memory responses  
186 appeared to have systematically different tuning? In particular, we sought to understand whether differences in  
187 memory responses could be explained by three types of noise: 1) reduced fMRI signal-to-noise; 2) retrieval failure  
188 on a subset of trials; 3) memory error. If perception and memory have the same fundamental response properties,  
189 but memory is subject to more noise, then adding noise to the perception data should yield responses that look  
190 like what we observed during memory. Thus, we started with perception data (mean and variance of each voxel's  
191 activity during perception) and tested whether we could generate responses that looked like memory data by adding  
192 one of the three types of noise. To simulate reduced fMRI signal-to-noise, we introduced additive noise to each  
193 voxel's perception response (Fig. 5a, left; see Methods). To simulate retrieval failure, we created some trials where  
194 the mean response was zero (Fig. 5b, left). To simulate memory error, we added angular noise to the peak location  
195 of the perception responses (Fig. 5c, left). For each of these types of simulation, we considered multiple levels of  
196 noise. To assess the simulation results, we analyzed all simulated datasets with the same procedures used for the  
197 real data and then counted the proportion of times the von Mises parameters derived from a simulation fell within  
198 the 95% confidence interval of the actual memory data (Fig. 5a-c, right).

199 First, using bootstrapped parameter estimates, we confirmed that the estimated signal-to-noise ratio (SNR)  
200 for perception parameter estimates was higher than for memory parameter estimates in every ROI. Perception  
201 SNR was between 1.3 and 1.6 times higher than memory SNR in each ROI, and between 2.2 and 4.3 times higher  
202 in vertices closest to the stimulus location. Given this, we simulated new perception data that precisely matched  
203 the empirical memory SNR for every surface vertex. We also simulated data with even lower SNR (higher noise)  
204 than what we observed during memory. As expected, simulating perception data with reduced SNR increased  
205 variance in the location, amplitude, and FWHM of the von Mises fits (Supplementary Fig. 1a). However, no level of  
206 SNR produced response profiles that matched the memory data well. In V1—the region where we observed the  
207 largest difference in FWHM between perception and memory—0% of the FWHM parameters in the memory SNR  
208 simulation approximated the actual memory data (Fig. 5a, right). In the noisiest simulation we performed (1/8 of  
209 the memory SNR), this figure was still only 10% (Fig. 5a, right). Similar results occurred for V2 and V3. These  
210 simulations demonstrate that low SNR cannot explain the pattern of memory responses we observed in early visual  
211 cortex.

212 Our SNR simulations also demonstrate that there are fundamental tradeoffs between capturing memory FWHM  
213 in early visual cortex and in capturing other aspects of the data. First, at high levels of noise (low levels of SNR),  
214 any modest increase in ability to capture V1-V3 FWHM was accompanied by a *decrease* in ability to capture  
215 FWHM in later visual areas (Fig. 5a, right). In these regions, FWHM was already equivalent during perception and  
216 memory, and artificially adding noise to the perception data harms this equivalence. Second, high noise simulations  
217 generated more noise in the location parameters than was actually observed in the memory data (Fig. 5a, right),  
218 resulting in unreliable location parameters in all ROIs.

219 We observed a similar pattern of results in the retrieval failure simulations. Very high rates of retrieval failure



**Figure 5. Differences between perception and memory are not explained by noise** (a) Left: We simulated the effect of low SNR by introducing additive noise to our perception data and asked whether this was sufficient to produce responses similar to what we observed during memory. Right: Proportion of simulations that produce location and FWHM parameters within the 95% confidence intervals of the memory data are plotted for decreasing signal-to-noise ratios (SNR) and for each ROI. SNR values ranged from the empirical SNR of the perception data (p), the empirical SNR of the memory data (m), or 1/2, 1/4, or 1/8 of the empirical SNR of the memory data. Even at extremely high noise levels, very few simulations generate FWHM parameters within the confidence intervals of the memory data in V1–V3. (b) Left: We simulated the effect of failing to perform the retrieval task by generating a perception dataset where a subset of trials had a mean BOLD response of zero. Right: Data are plotted as in a, with increasing large numbers of failed trials on the x axis. As in a, even at extremely high rates of failed retrieval, FWHM parameters in V1–V3 rarely fall within the memory confidence interval (c) Left: We simulated the effect of memory error by adding angular noise to the peak location of the perception responses. Right: Data are plotted in a, with increasing large amounts of memory error on the x axis. Implausibly large amounts of memory error are needed to generate FWHM parameters that fall within the memory confidence intervals  $\geq 50\%$  of the time in V1–V3. In all panels, increased noise produced a worse match to memory FWHM in hV4, LO and V3ab, as well as unreliable location parameters in all ROIs.

were required to generate any FWHM parameters that were sufficiently wide to match the memory data in V1 (Supplementary Fig. 1b). Only when simulating retrieval failure in  $>=50\%$  of all trials, did this number exceed 0% (Fig. 5b, right). Similar to the SNR simulation, any improvement in ability to capture the V1 FWHM data with increased retrieval failure was offset by a decline in ability to capture FWHM in late visual areas (Fig. 5b, right), where responses became much wider than what was observed empirically during memory. Again, as in the SNR simulation, high rates of retrieval failure were also associated with location parameters that were far noisier than what we observed (Fig. 5b, right). Thus, subjects experiencing retrieval failure on a subset of trials does not explain our memory data.

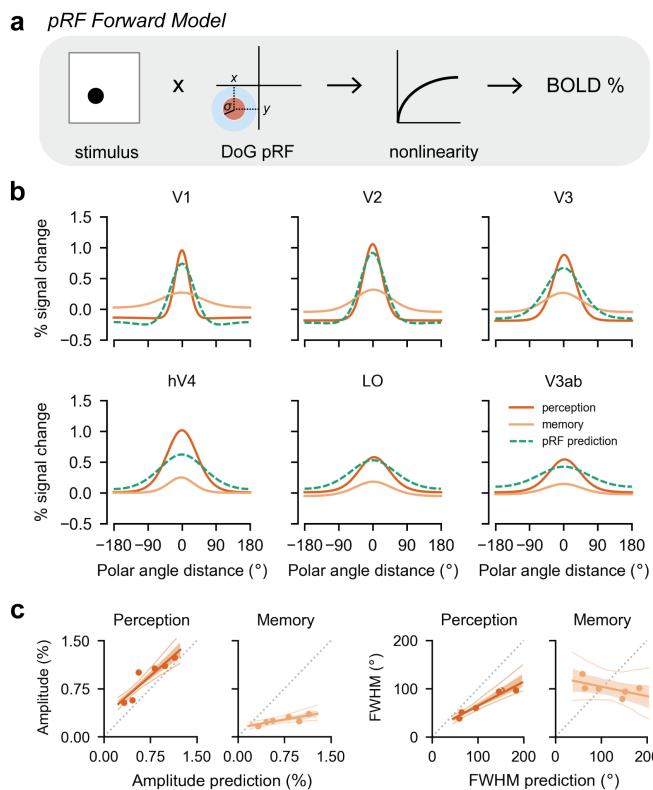
Finally, we considered the memory error simulation. Compared to the other simulations, this simulation produced a better match to memory FWHM in V1 when assuming high levels of noise (Supplementary Fig. 1c). Still, in the best performing simulation only 57% of the V1 FWHM parameters approximated the memory data (Fig. 5c, right). Critically, the magnitude of memory error in this simulation was implausibly high. The standard deviation of memory errors around the true value was  $45^\circ$ , meaning that simulated memories were within the correct quadrant only 68% of time. Given that subjects were trained to discriminate remembered locations up to  $15^\circ$  (see Methods), errors of this magnitude and frequency are exceedingly unlikely. At a more plausible  $15^\circ$  standard deviation of memory error, 0% of simulations approximated the memory data (Fig. 5c, right). Further, similar to the other simulations, improvements in the ability to capture V1 FWHM with high levels of angular error were accompanied by decreases

in the ability to capture FWHM in later areas and by decreases in location parameter reliability beyond what we observed empirically (Fig. 5c, right). Thus, subjects experiencing a small, variable amount of memory error does explain our memory data.

Collectively, these simulations demonstrate that our results are unlikely to be caused by a simple source of measurement noise (reduced SNR) or cognitive noise (failed retrieval, memory error). In each of the three simulations, the amount of noise required to make even modest gains in our ability to account for the V1 memory FWHM was implausibly large. Further, in all three cases, increases in the ability to account for V1 FWHM were accompanied by decreases in the ability to account for FWHM in higher visual areas and to recover location parameters that were as reliable as our actual data.

#### pRF models accurately predict perception but not memory responses

Next, we evaluated how well perception and memory responses matched the predictions of a pRF model. To do this, we used a modified version of each subject's pRF model to generate predicted cortical responses to each of the four experimental stimuli (Fig. 6a; see Methods). The pRF model we used to generate predictions is a novel variant of existing pRF models: we added a Difference of Gaussian pRF shape (Zuiderbaan et al., 2012) with a fixed positive to negative Gaussian size ratio (1:2) and amplitude ratio (2:1) to our solved nonlinear compressive spatial summation (CSS) model (Kay et al., 2013b). The predictions from the model were analyzed with the same procedure as the data, yielding von Mises fits to the predicted data (Fig. 6b). Model predictions from simpler pRF models are shown in Supplementary Figure 2.



**Figure 6. pRF forward model captures perception but not memory responses.** (a) We used our pRF model to generate the predicted BOLD response to each of our experimental stimuli. The model assumes a Difference of Gaussians pRF shape, with a fixed positive to negative Gaussian size ratio (1:2) and amplitude ratio (2:1). The model also incorporates a compressive nonlinearity. (b) Predicted polar angle response functions are plotted for the pRF model (green dashed lines), alongside the functions fit to the perception and memory data (dark and light orange, reproduced from Fig. 4b). The model predictions are closer to the perception data than the memory data in all visual areas. (c) Predicted versus observed amplitude (left) and FWHM (right), plotted separately for perception and memory. Each dot represents an ROI. The shaded region is the 68% CI from bootstrapping linear fits across participants, and the thin lines indicate the 95% confidence intervals. For both the amplitude and FWHM, the perception data lie relatively close to the pRF model predictions (dashed grey lines), whereas the memory data do not.

Qualitatively, the pRF model predictions agree with the perception data but not the memory data (Fig. 6b). Several specific features of the perception data are well captured by the model. First, the model predicts the highest amplitude response at cortical sites with pRFs near the stimulus location (peak at 0°). Second, the model predicts increasingly wide response profiles from the early to late visual areas. Third, it predicts higher amplitudes in early compared to late areas. Finally, the model predicts negative responses in the surround locations of V1-V3 but not higher visual areas. This is particularly interesting given that all voxel pRFs were implemented with a negative surround of the same size and amplitude relative to the center Gaussian. This suggests that voxel-level parameters and population-level responses can diverge (Sprague & Serences, 2013, see also). Though not the focus of this analysis, we note that the model predictions are not perfect. The model predicts slightly lower amplitudes and larger FWHM than is observed in the perception data (Fig. 6b). These discrepancies may be due to differences between the stimuli used in the main experiment and those used in the pRF experiment or to differences in the task (attending fixation during the pRF experiment vs attending the stimulus during the main experiment).

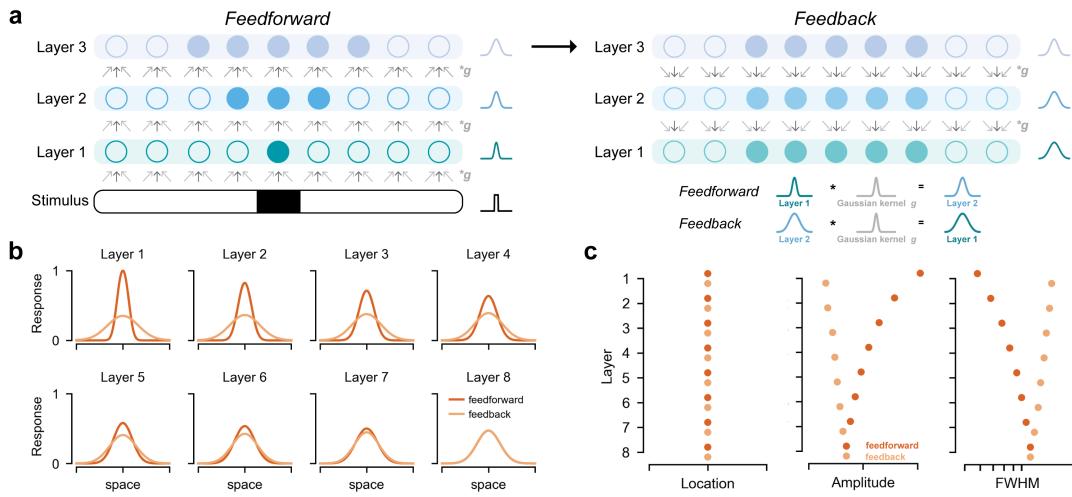
Critically, the model accurately captures the properties of memory responses that are shared with perception responses (the peak location), but not the distinct properties (Fig. 6b). These failures are especially clear when comparing the predicted amplitude and FWHM from the pRF model with the observed amplitudes and FWHMs for perception and memory. While there is a positive slope between the predicted amplitude and both the perception ( $\beta = 0.84$ , 95% CI: [0.56, 1.15]) and memory amplitudes ( $\beta = 0.17$ , 95% CI: [0.056, 0.32]), the slopes differ substantially (Fig. 6c). The perception amplitudes have a slope closer to 1, indicating good agreement with the model predictions, while the memory data have a slope closer to 0, indicating weak agreement. Similarly, the predicted FWHM is strongly and positively related to the perception FWHM ( $\beta = 0.50$ , 95% CI: [0.36, 0.76]), but weakly and negatively related to the memory FWHM ( $\beta = -0.20$ , 95% CI: [-0.67, 0.26]; Fig. 6c). These analyses strongly support our interpretation of the data in Figure 4b,c to mean that memory and perception have distinct spatial tuning properties. The critical advantage of using pRF models is that they explicitly incorporate known properties of feedforward spatial processing in visual cortex. Because our pRF model fails to account for the memory responses we observed, we can conclude that memory reactivation violates the assumptions of feedforward processes that well characterize perceptual activation. A plausible explanation for this failure is that memory retrieval involves a fundamentally different origin and cascade of information through visual cortex, a possibility we explore in detail in the next section.

## Perception and memory responses can be simulated with a bidirectional hierarchical model

Cortical activity during perception arises from a primarily feedforward process that originates with the retina and that accumulates additional spatial pooling in each cortical area, resulting in increasingly large receptive fields (Gattass et al., 2005; Wandell & Winawer, 2015). In contrast, memory reinstatement is hypothesized to begin with the hippocampus (Marr, 1971; O'Reilly & McClelland, 1994), a region bidirectionally connected to high-level visual areas in ventral temporal cortex via the medial temporal lobe cortex (Van Hoesen & Pandya, 1975; Felleman & Essen, 1991; Suzuki & Amaral, 1994). Reinstated cortical activity is then thought to propagate backwards through visual cortex (Naya et al., 2001; Linde-Domingo et al., 2019; Dijkstra et al., 2019; Hindy et al., 2016). Here, we explored whether a simple hierarchical model with spatial pooling could be adapted to account for both our perception and memory results by manipulating the direction of information flow.

We first constructed a linear feedforward hierarchical model of spatial processing in neocortex. In this model, the activity in each layer was created by convolving the activity from the previous layer with a fixed Gaussian kernel (Fig. 7a; see Methods). Beginning with a boxcar stimulus, we cascaded this convolutional operation to simulate 8 layers of the network (Fig. 7b). In this simple demonstration, the size of the convolutional kernel was fixed, not fit to the data. Nonetheless, the pattern of feedforward responses qualitatively matches our fMRI observations during perception. The location of the peak response is unchanged across layers, but response functions become wider and lower in amplitude in higher layers (Fig. 7b,c)—precisely as we observed in our actual data (Fig. 4b,c).

We then explored whether backwards propagation of reinstated activity in our hierarchical model could account for our memory data. To do this, we assumed that feedforward and feedback connections in the model were reciprocal, meaning that the convolutional kernel was the same in feedforward and feedback direction. We assumed perfect reinstatement in the top layer, and thus began the feedback simulation by duplicating the feedforward activity



**Figure 7. Perception and memory responses can be simulated with a bidirectional hierarchical model.** (a) Illustration of stimulus-driven activity propagating through a linear hierarchical network model in the feedforward direction (left) and mnemonic activity propagating through the model in the feedback direction (right). In both cases, a given layer's activity is generated by convolving the previously active layer's activity with a fixed Gaussian kernel. The feedforward simulation began with a boxcar stimulus. The feedback simulation began with duplication of the feedforward activity from the final layer. (b) Results from feedforward and feedback simulations in an 8 layer network, plotted in the conventions of Figure 4b. The feedforward simulation parallels our observations during perception, and the feedback simulation parallels our observations during memory. (c) Location, amplitude, and FWHM parameters for each layer, plotted separately for feedforward and feedback simulations. Location is preserved across layers in the feedforward and feedback direction. Note that FWHM become progressively wider in later layers in the feedforward direction and in earlier layers in the feedback direction. This results in large differences in FWHM between feedforward and feedback activity in early layers. These trends closely follow our observations in Figure 4c.

from the final layer. Starting with this final layer activity, we convolved each layer's activity with the same Gaussian kernel to generate earlier layers' activity (Fig. 7a). The properties of the simulated activity (Fig. 7b,c) bear a striking resemblance to those of the observed memory data (Fig. 4b,c). First, simulated feedback activity had a preserved peak location across layers (Fig. 7c, left), similar to the memory data. Second, simulated feedback activity was wider and lower amplitude than feedforward activity overall (Fig. 7c, middle and right)—just as the memory data had wider and lower amplitude responses than the perception data. Third, the increase in FWHM across layers was smaller in the feedback direction than in the feedforward direction, and it reversed direction with respect to the visual hierarchy (Fig. 7c, right). This small effect of reversal is particularly interesting given that this trend was numerically present in our memory data but not statistically reliable at our sample size. Finally, the difference between feedforward and feedback FWHM was maximal in the earliest layers (Fig. 7c, right), just as the difference between our perception and memory data was maximal in V1. This simulation suggests that the distinct spatial profile of mnemonic responses in visual cortex may be a straightforward consequence of reversing the flow of information in a system with hierarchical structure and reciprocal connectivity, and that spatial pooling accumulated during feedforward processing may not be inverted during reinstatement. More broadly, these results demonstrate that models of the visual system may be useful for probing the mechanisms that support and constrain visual memory.

319 **Discussion**

320 In the current work, we combined empirical and modeling approaches to explore how long-term spatial memories  
321 are represented in the human visual system. By using computational models of spatial encoding to compare  
322 perceptual and mnemonic BOLD activity, we provide strong evidence that visual memory, like visual perception,  
323 produces retinotopically-mapped activation throughout visual cortex. Critically, however, we also identified systematic  
324 differences in the population spatial tuning properties of perceptual and mnemonic activity. Compared to perceptual  
325 responses, mnemonic responses were lower in amplitude in all visual areas. Further, while we observed a three-  
326 fold change in spatial precision from early to late visual areas during perception, mnemonic responses violated  
327 this pattern. Instead, mnemonic responses displayed consistent spatial precision across visual areas. Notably,  
328 simulations showed that neither reduced SNR, nor failure to retrieve on some trials, nor memory error could  
329 account for this difference. We speculate, instead, that this difference arises from a reversal of information flow in a  
330 hierarchically organized and reciprocally connected visual cortex. To support this, we show that top-down activation  
331 in a simple hierarchical model elicits a systematically different pattern of responses than bottom-up activation. These  
332 simulations reproduce the properties we observe during both perception and memory. Together, these results reveal  
333 novel properties of memory-driven activity in visual cortex that suggest specific computational processes governing  
334 visual cortical responses during memory retrieval.

355 **Advantages of using encoding models to parameterize memory representations**

356 Much work in neuroscience has been dedicated to the question of how internally-generated stimulus representations  
357 are coded in the brain. Early neuroimaging work established that sensory cortices are recruited during imagery and  
358 memory tasks (Kosslyn et al., 1995; O’Craven & Kanwisher, 2000; Wheeler et al., 2000), moving the field away  
359 from purely symbolic accounts of memory (e.g. Pylyshyn, 2002). More recently, memory researchers have favored  
360 decoding and pattern similarity approaches over univariate activation analyses to examine the content of retrieved  
361 memories (Polyn et al., 2005; Kuhl et al., 2011; Favila et al., 2018). While these approaches are powerful, they do  
362 not explicitly specify the form mnemonic activity should take, and many activation schemes can lead to successful  
363 decoding or changes in pattern correlations. In the present work, we leveraged encoding models from visual  
364 neuroscience, specifically stimulus-referred pRF models, to examine and account for memory-triggered activity in  
365 visual cortex. In contrast to decoding or pattern similarity approaches, encoding models predict the activity evoked  
366 in single voxels in response to sensory or cognitive manipulations using a set of explicit mathematical operations  
367 (Naselaris et al., 2011). Spatial encoding models have proved particularly powerful because space is coded in the  
368 human brain at a scale that is well-matched to the millimeter sampling resolution of fMRI (Engel et al., 1994; Sereno  
369 et al., 1995; Dougherty et al., 2003). Despite the power of such encoding models, relatively little work has applied  
370 these models to questions about long-term memory (c.f. Thirion et al., 2006; Naselaris et al., 2015; Breedlove et al.,  
371 2018). Here, using this approach, we revealed novel properties of memory responses in visual cortex that decoding  
372 approaches have missed. Most notably, we found that memory activity was characterized by a different pattern  
373 of spatial precision across regions than perceptual activity. Because spatial parameters such as polar angle are  
374 explicitly modeled in pRF models, we were able to quantify and interpret these differences.

355 Our results have important implications for the study of memory reactivation. First, our findings suggest that the  
356 specific architecture of a sensory system may constrain what memory reactivation looks like in that system. Though  
357 memory reactivation is often studied in sensory domains, the architecture of these systems is not usually considered  
358 when interpreting reactivation effects. Here, we propose that hierarchical spatial pooling in visual cortex produces a  
359 systematic and distinct pattern of memory reactivation that cannot be attributed to retrieval failure or memory error.  
360 However, whether this architecture has any consequences for memory behavior is not clear from the present study.  
361 This question will be critical for future studies to address. Second, our results advocate for shifting away from the  
362 concept of memory reactivation as it has been understood and applied in the field of neuroimaging. Most previous  
363 work has focused on identifying *similarities* between the neural substrates of visual perception and visual memory.  
364 These studies have been successful in that they have produced many positive findings of memory reactivation  
365 in human visual cortex (Kosslyn et al., 1995; O’Craven & Kanwisher, 2000; Wheeler et al., 2000; Slotnick et al.,  
366 2005; Polyn et al., 2005; Kuhl et al., 2011; Bosch et al., 2014; Waldhauser et al., 2016; Lee et al., 2018; Bone

367 et al., 2018). However, much of this work implicitly assumes that any mismatch between perception and memory  
368 is due to the fact that memory reactivation is either inherently low fidelity or susceptible to noise (Pearson et al.,  
369 2015), or is a subset of the perceptual response (Wheeler et al., 2000). Our results demonstrate that, at least  
370 in the spatial domain, this is not the case, and that systematic differences beyond noise exist. These results are  
371 broadly consistent with other recent findings suggesting computational differences between perception and memory  
372 derived from behavior (Bloem et al., 2018) and multivoxel pattern differences in perceived and remembered object  
373 representations measured with fMRI (Lee et al., 2012). Ultimately, the field should strive to identify, quantify, and  
374 explain these differences in order to fully understand the neural basis of memory. Using encoding models borrowed  
375 from sensory neuroscience to parameterize the differences between perception and memory may prove a fruitful  
376 way of making progress on this goal.

### 377 Why do perceptual and mnemonic representations differ in visual cortex?

378 Despite the usefulness of encoding models like pRF models for quantifying neural responses in a stimulus-referred  
379 space, these models may not provide a natural explanation for *why* perception and memory responses differ.  
380 We show in Figure 6 that pRF models fail to capture the aspects of memory responses that are distinct from  
381 perceptual responses: namely, the dramatic change in spatial precision. While it would be possible to fit separate  
382 pRF parameters to memory data to improve the ability of the model to accurately predict memory responses, this still  
383 would not explain why these parameters or responses differ. How then can we account for this? We were particularly  
384 intrigued by the possibility that differences between memory and perception activity are a direct consequence of the  
385 direction of processing in hierarchically-organized cortex. Hierarchical structure and feedback processing are not  
386 typically directly simulated in a pRF model but there is considerable evidence to suggest these factors are of interest.  
387 Studies of anatomical connectivity provide evidence that the visual system is organized approximately hierarchically  
388 (Felleman & Essen, 1991; Barone et al., 2000; for other perspectives see Zeki, 2015; Hilgetag & Goulas, 2020),  
389 and that most connections within the visual system are reciprocal (Felleman & Essen, 1991). Studies also show  
390 that the hippocampus sits atop the highest stage of the visual hierarchy, with reciprocal connections to high-level  
391 visual regions via the medial temporal lobe cortex (Van Hoesen & Pandya, 1975; Felleman & Essen, 1991; Suzuki  
392 & Amaral, 1994). These observations make the prediction that initial drive from the hippocampus during memory  
393 retrieval should propagate backwards through the visual system. Neural recordings from the macaque (Naya et al.,  
394 2001) and human (Hindry et al., 2016; Linde-Domingo et al., 2019; Dijkstra et al., 2019), as well as computational  
395 modeling (Horikawa & Kamitani, 2017) support this idea.

396 Based on these observations and our hypothesis, we constructed a hierarchical network model in which we could  
397 simulate top-down activity. Though this model shares some features of hierarchical models of object recognition  
398 (Riesenhuber & Poggio, 1999; Serre et al., 2007), we emphasize that it is much simpler. Our model is entirely  
399 linear, its parameters are fixed *a priori* (not the result of training), and it encodes only one stimulus feature: space  
400 (Kay et al., 2013b). Critically, in contrast to pRF models, which express each region's activity as a function of the  
401 stimulus, our model expresses each region's activity as a function of the previous region's activity (Fukushima, 1980;  
402 Riesenhuber & Poggio, 1999), and can therefore simulate both feedforward and feedback processes (Heeger, 2017).  
403 While highly simplified, the simulations we performed in this network captured the dominant features of our data,  
404 providing a parsimonious explanation for our observations. Interestingly, our simulations also indicate that some  
405 trends present in our data warrant further investigation. For instance, while we could not conclude that the earliest  
406 visual cortical areas had the least precise responses during memory (a *reversal* of the perception pattern), our  
407 simulations suggest that this effect should be present, albeit significantly weaker than in the feedforward direction.  
408 Future work should target this small effect with a sufficiently powered experiment.

409 Our simulations also raise interesting questions and predictions about the consequences of visual cortical  
410 architecture for cognition. First, why have a hierarchical architecture in which the detailed information present in  
411 early layers cannot be reactivated? The hierarchical organization of the visual system is thought to give rise to the  
412 low-level feature invariance required for object recognition (Riesenhuber & Poggio, 1999; Serre et al., 2007). Our  
413 results raise the possibility that the benefits of such an architecture for recognition outweigh the cost of reduced  
414 precision in top-down responses. Whether the extent of this tradeoff differs between healthy individuals or between  
415 healthy and neuropsychiatric populations, and what consequences this structure has for behavior, are interesting

416 questions for future research. Second, how is it that humans have spatially precise memories if visual cortical  
417 responses do not reflect this? One possibility is that read-out mechanisms are not sensitive to all of the properties  
418 of mnemonic activity we measured. For instance, memory decisions could be driven exclusively by the neural  
419 population with the strongest response (e.g. those at the peak of the polar angle response functions). Another  
420 possibility is that regions without hierarchical structure do not exhibit these properties and reactivation in these  
421 other regions is preferentially used to guide memory-based behavior. These, and other possibilities should be  
422 directly explored in future work. Finally, our hierarchical simulations highlight the need to carefully separate the  
423 contribution of visual cortical architecture on reactivation from the effects of cognitive manipulations or effects  
424 occurring upstream of visual cortex (e.g. in the hippocampus).

## 425 Relation to other forms of memory and attention

426 Sensory reactivation during long-term memory retrieval has parallels to sensory engagement in other forms of  
427 memory such as iconic memory and working memory. Nonetheless there may also be differences in the specific  
428 way that sensory circuits are used across these forms of memory. One critical factor may be how recently the  
429 sensory circuit was activated by a stimulus at the time of memory retrieval. In iconic memory studies, very detailed  
430 information can be retrieved if probed within a second of the sensory input (Sperling, 1960). In working memory  
431 studies, sensory activity is thought to be maintained by active mechanisms from stimulus encoding through a  
432 seconds-long delay. Using similar methods to the ones we use here (Sprague & Serences, 2013; Ester et al., 2013),  
433 many working memory studies have shown that early visual areas contain retinotopically specific signals throughout  
434 a delay period (Sprague & Serences, 2013; Sprague et al., 2014; Rahmati et al., 2017), paralleling our findings. In  
435 imagery studies, eye-specific circuits presumed to be in V1 can be re-engaged if there is a delay of 5 minutes or  
436 less from when the subject viewed stimuli through the same eye, but not if there is a delay of 10 minutes (Ishai  
437 & Sagi, 1995). Hippocampally-dependent memory retrieval is thought to be capable of engaging visual cortex at  
438 much longer delays. Given that the mechanism for engaging sensory cortex may differ across these different forms  
439 of memory, the question of how similar sensory activation is across these timescales remains an important open  
440 question. For example, shorter-term forms of memory might, in principle, cause more spatially-specific reactivation  
441 in early visual cortex than what we observed in long-term memory. Informal comparisons between our data and  
442 stimulus reconstructions made from working memory delay period activity (Sprague & Serences, 2013; Rahmati  
443 et al., 2017; Rademaker et al., 2019) suggest this may be the case, but a direct comparison is warranted. The  
444 current study offers a quantitative approach for directly comparing spatial tuning properties across different cognitive  
445 processes, and could be extended to include multiple forms of memory within the same experiment.

446 Are the spatial responses we observed during memory retrieval better characterized as long-term memory  
447 reactivation or as a special case of (memory-guided) spatial attention? Our results raise interesting questions  
448 about whether long-term spatial memory and endogenous spatial attention share mechanisms for modulating the  
449 response of visual cortical populations. In typical endogenous spatial attention tasks, subjects are explicitly cued to  
450 the most likely location of an upcoming stimulus prior to being presented with a difficult visual judgment (Carrasco,  
451 2011). fMRI studies have repeatedly found that spatial attention enhances visually-evoked responses in visual  
452 cortex (Somers et al., 1999; Gandhi et al., 1999; Buracas & Boynton, 2007; Li et al., 2008). Similar to our results,  
453 spatial attention has also been shown to elicit spatially localized activation in the absence of any visual stimulation  
454 (Luck et al., 1997; Kastner et al., 1999; Chawla et al., 1999; Ress et al., 2000). It is at least logically possible for  
455 attention and memory to dissociate. Most endogenous attention tasks have no memory component since the cue  
456 explicitly represents the attended location. In contrast, in most episodic memory tasks the association between a  
457 cue and a stimulus is intentionally arbitrary so that it must be acquired and retrieved in a hippocampally-dependent  
458 manner. However, it is possible that spatial attention and memory processes only differ in their dependency on  
459 the hippocampus to retrieve the target location. Once this target location is determined, the same mechanisms  
460 could be used to initiate enhanced processing of the target location in sensory areas. Future experiments and  
461 modeling efforts should determine whether memory-driven and attention-driven activations in visual cortex differ,  
462 and whether it's possible to develop a model of top-down processing in visual cortex that can account for both sets  
463 of observations.

464 **Conclusion**

465 In the current work, we provide novel empirical evidence that memory retrieval elicits systematically different  
466 activation in human visual cortex compared to visual perception. Using simulations and a network model of cortex,  
467 we argue that these distinctions arise from a reversal of information flow within a hierarchically structured visual  
468 system. Collectively, this work makes progress on providing a detailed account of reactivation in visual cortex and  
469 sheds light on the broader computational principles that guide top-down processes in sensory systems.

470 **Methods**

471 **Subjects**

472 Nine human subjects participated in the experiment (5 males, 22–46 years old). All subjects had normal or  
473 correct-to-normal visual acuity, normal color vision, and no MRI contraindications. Subjects were recruited from  
474 the New York University community and included author S.E.F and author J.W. All subjects gave written informed  
475 consent to procedures approved by the New York University Institutional Review Board prior to participation. No  
476 subjects were excluded from data analysis.

477 **Stimuli**

478 Experimental stimuli included nine unique radial frequency patterns (Fig. 1a). We first generated patterns that  
479 differed along two dimensions: radial frequency and amplitude. We chose stimuli that tiled a one dimensional  
480 subspace of this two dimensional space, with radial frequency inversely proportional to amplitude. The nine chosen  
481 stimuli took radial frequency and amplitude values of: [2, .9], [3, .8], [4, .7], [5, .6], [6, .5], [7, .4], [8, .3], [9, .2], [10,  
482 .1]. We selected four of these stimuli to train subjects on in the behavioral training session and to appear in the  
483 fMRI session. For every subject, those stimuli were: [3, .8], [5, .6], [7, .4], [9, .2]; (radial frequency, amplitude). The  
484 remaining five stimuli were used as lures in the test trials of the behavioral training session. Stimuli were saved as  
485 images and cropped to the same size.

486 **Experimental procedure**

487 The experiment began with a behavioral training session, during which subjects learned four paired associates (Fig.  
488 1). Specifically, subjects learned that four colored fixation dot cues were uniquely associated with four spatially  
489 localized radial frequency patterns. An fMRI session immediately followed completion of the behavioral session  
490 (Fig. 2a). During the scan, subjects participated in two types of functional runs (approximately 3.5 min each): (1)  
491 perception, where they viewed the cues and associated spatial stimuli; and (2) memory, where they viewed only the  
492 fixation cues and recalled the associated spatial stimuli. Details for each of these phases are described below. A  
493 separate retinotopic mapping session was also performed for each subject (Fig. 2b), which is described in the next  
494 section.

495 **Behavioral training**

496 For each subject, the four radial frequency patterns were first randomly assigned to one of four polar angle locations  
497 in the visual field ( $45^\circ$ ,  $135^\circ$ ,  $225^\circ$ , or  $315^\circ$ ) and to one of four colored cues (orange, magenta, blue, green; Fig. 1b).  
498 Immediately before scanning, subjects learned the association between the four colored cues and the four spatially  
499 localized stimuli through interleaved study and test blocks (Fig. 1c). Subjects alternated between study and test  
500 blocks, completing a minimum of four blocks of each type. Subjects were required to reach at least 95% accuracy,  
501 and performed additional rounds of study-test if they did not reach this threshold after four test blocks.

502 During study blocks, subjects were presented with the associations. Subjects were instructed to maintain central  
503 fixation and to learn each of the four associations in anticipation of a memory test. At the start of each study trial  
504 (Fig. 1c), a central white fixation dot (radius = 0.1 dva) switched to one of the four cue colors. After a 1 sec delay,  
505 the associated radial frequency pattern appeared at  $2^\circ$  of eccentricity and its assigned polar angle location in the

506 visual field. Each pattern image subtended 1.5 dva and was presented for 2 sec. The fixation dot then returned  
507 to white, and the next trial began after a 2 sec interval. No subject responses were required. Each study block  
508 contained 16 trials (4 trials per association), presented in random order.

509 During test blocks, subjects were presented with the colored fixation dot cues and tested on their memory for the  
510 associated stimulus pattern and spatial location. Subjects were instructed to maintain central fixation and to try to  
511 covertly recall each stimulus when cued, and then to respond to the test probe when prompted. At the start of each  
512 test trial (Fig. 1c), the central white fixation dot switched to one of the four cue colors. This cue remained on the  
513 screen for 2.5 sec while subjects attempted to covertly retrieve the associated stimulus. At the end of this period,  
514 a test stimulus was presented at 2° of eccentricity for 2 sec. Then, subjects were cued to make two consecutive  
515 responses to the test stimulus: whether it was the correct radial frequency pattern (yes/no) and whether it was  
516 presented at the correct polar angle location (yes/no). Each test stimulus had a 50% probability of being the correct  
517 pattern. Incorrect patterns were drawn randomly from the three patterns associated with other cues and the five lure  
518 patterns (Fig. 1a). Each test stimulus had a 50% probability of being in the correct polar angle location, which was  
519 independent from the probability of being the correct pattern. Incorrect polar angle locations were drawn from the  
520 three locations assigned to the other patterns and 20 other evenly spaced locations around the visual field (Fig. 1a).  
521 This placed the closest spatial lure at 15° of polar angle away from the correct location. Responses were solicited  
522 from the subject with the words "Correct pattern?" or "Correct location?" displayed centrally in white text. The order  
523 of these queries was counterbalanced across test blocks. Subjects responses were recorded on a keyboard with  
524 a maximum response window of 2 sec. Immediately after a response was made or the response window closed,  
525 the color of the text turned black to indicate an incorrect response if one was made. After this occurred for both  
526 queries, subjects were presented with the colored fixation dot cue and correct spatially localized pattern for 1 sec as  
527 feedback. This feedback occurred for every trial, regardless of subject responses to the probe. Each test block  
528 contained 16 trials (4 trials per association), presented in random order.

#### 529 **fMRI session**

530 During the fMRI session, subjects participated in two types of functional runs: perception and memory retrieval (Fig.  
531 2a). Subjects completed 5–6 runs each of perception and memory in an interleaved order. This amounted to 40–48  
532 repetitions of perceiving each stimulus and of remembering each stimulus per subject.

533 During perception runs, subjects viewed the colored fixation dot cues and the radial frequency patterns in their  
534 learned locations. Subjects were instructed to maintain central fixation and to perform a one-back task on the  
535 stimuli. The purpose of the one-back task was to encourage covert stimulus-directed attention on each trial. At  
536 the start of each perception trial (Fig. 2a, top), a central white fixation dot (radius = 0.1 dva) switched to one of  
537 the four cue colors. After a 0.5 sec delay, the associated radial frequency pattern appeared at 2° of eccentricity  
538 and its assigned polar angle location in the visual field. Each pattern subtended 1.5 dva and was presented for 2.5  
539 sec. The fixation dot then returned to white and the next trial began after a variable interval. Intervals were drawn  
540 from an approximately geometric distribution sampled at 3, 4, 5, and 6 sec with probabilities of 0.5625, 0.25, 0.125,  
541 and 0.0625 respectively. Subjects indicated when a stimulus repeated from the previous trial using a button box.  
542 Responses were accepted during the stimulus presentation or during the interstimulus interval. Each perception run  
543 contained 32 trials (8 trials per stimulus). The trial order was randomized for each run, separately for every subject.

544 During memory runs, subjects viewed the colored fixation dot cues and recalled the associated patterns in  
545 their learned spatial locations. Subjects were instructed to maintain central fixation, to use the cues to initiate  
546 recollection, and to make a subjective judgment about the vividness of their memory on each trial. The purpose of  
547 the vividness task was to enforce attention to the remembered stimulus on each trial. At the start of each memory  
548 trial (Fig. 2a, top), the central white fixation dot switched to one of the four cue colors. This cue remained on the  
549 screen for a recollection period of 3 sec. The fixation dot then returned to white and the next trial began after a  
550 variable interval. Subjects indicated whether the stimulus associated with the cue was vividly remembered, weakly  
551 remembered, or not remembered using a button box. Responses were accepted during the cue presentation or  
552 during the interstimulus interval. Each memory run contained 32 trials (8 trials per stimulus). For a given subject,  
553 each memory run's trial order and trial onsets were exactly matched to one of the perception runs. The order of  
554 these matched memory runs was scrambled relative to the order of the perception runs.

555 **Retinotopic mapping procedure**

556 Each subject completed either 6 or 12 identical retinotopic mapping runs in a separate fMRI session from the  
557 main experiment (Fig. 2b, top). Stimuli and procedures for the retinotopic mapping session were based on those  
558 used by the Human Connectome Project (Benson et al., 2018) and were identical to those reported in Benson &  
559 Winawer (2018). During each functional run, bar apertures on a uniform gray background swept across the central  
560 24 degrees of the subject's visual field (circular aperture with a radius of 12 dva). Bar apertures were a constant  
561 width (1.5 dva) at all eccentricities. Each sweep began at one of eight equally spaced positions around the edge  
562 of the circular aperture, oriented perpendicularly to the direction of the sweep. Horizontal and vertical sweeps  
563 traversed the entire diameter of the circular aperture while diagonal sweeps stopped halfway and were followed by  
564 a blank period. A full-field sweep or half-field sweep plus blank period took 24 s to complete. One functional run  
565 contained 8 sweeps, taking 192 s in total. Bar apertures contained a grayscale pink noise background with randomly  
566 placed faces, scenes, objects, and words at a variety of sizes. Noise background and stimuli were updated at a  
567 frequency of 3 Hz. Each run of the task had an identical design. Subjects were instructed to maintain fixation on a  
568 central dot and to use a button box to report whenever the dot changed color. Color changes occurred on average  
569 every 3 s.

570 **MRI acquisition**

571 Images were acquired on a 3T Siemens Prisma MRI system at the Center for Brain Imaging at New York University.  
572 Functional images were acquired with a T2\*-weighted multiband EPI sequence with whole-brain coverage (repetition  
573 time = 1 s, echo time = 37 ms, flip angle = 68°, 66 slices, 2 x 2 x 2 mm voxels, multiband acceleration factor  
574 = 6, phase-encoding = posterior-anterior) and a Siemens 64-channel head/neck coil. Spin echo images with  
575 anterior-posterior and posterior-anterior phase-encoding were collected to estimate the susceptibility-induced  
576 distortion present in the functional EPIs. Between one and three whole-brain T1-weighted MPRAGE 3D anatomical  
577 volumes (.8 x .8 x .8 mm voxels) were also acquired for seven subjects. For two subjects, previously acquired  
578 MPRAGE volumes (1 x 1 x 1 mm voxels) from a 3T Siemens Allegra head-only MRI system were used.

579 **MRI processing**

580 **Preprocessing**

581 Anatomical and functional images were preprocessed using FSL v5.0.10 (Smith et al., 2004) and Freesurfer v5.3.0  
582 (Fischl, 2012) tools implemented in a Nipype workflow (Gorgolewski et al., 2011). To correct for head motion, each  
583 functional image acquired in a session was realigned to a single band reference image and then registered to  
584 the spin echo distortion scan acquired with the same phase encoding direction. The two spin echo images with  
585 reversed phase encoding were used to estimate the susceptibility-induced distortion present in the EPIs. For each  
586 EPI volume, this nonlinear unwarping function was concatenated with the previous spatial registrations and applied  
587 with a single interpolation. Freesurfer was used to perform segmentation and cortical surface reconstruction on  
588 each subject's average anatomical volume. Registration from the functional images to each subject's anatomical  
589 volume was performed using boundary-based registration. Preprocessed functional time series were then projected  
590 onto each subject's reconstructed cortical surface.

591 **GLM analyses**

592 Beginning with each subject's surface-based time series, we used GLMdenoise (Kay et al., 2013a) to estimate the  
593 neural pattern of activity evoked by perceiving and remembering every stimulus (Fig. 2a). GLMdenoise improves  
594 signal-to-noise ratios in GLM analyses by identifying a pool of noise voxels whose responses are unrelated to the  
595 task and regressing them out of the time series. This technique first converts all time series to percent signal change  
596 and determines an optimal hemodynamic response function for all vertices using an iterative linear fitting procedure.  
597 It then identifies noise vertices as vertices with negative  $R^2$  values in the task-based model. Then, it derives noise  
598 regressors from the noise pool time series using principal components analysis and iteratively projects them out of  
599 the time series of all vertices, one noise regressor at a time. The optimal number of noise regressors is determined

600 based on cross-validated  $R^2$  improvement for the task-based model. We estimated two models using this procedure.  
601 We constructed design matrices for the perception model to have four regressors of interest (one per stimulus),  
602 with events corresponding to stimulus presentation. Design matrices for the memory model were constructed the  
603 same way, with events corresponding to the cued retrieval period. These models returned parameter estimates  
604 reflecting the BOLD amplitude evoked by perceiving or remembering a given stimulus versus baseline for every  
605 vertex on a subject's cortical surface (Fig. 2a, bottom).

606 **Fitting pRF models**

607 Images from the retinotopic mapping session were preprocessed as above, but omitting the final step of projecting  
608 the time series to the cortical surface. Using these time series, nonlinear symmetric 2D Gaussian population  
609 receptive field (pRF) models were estimated in Vistasoft (Fig. 2b), as described previously (Dumoulin & Wandell,  
610 2008; Kay et al., 2013b). We refer to this nonlinear version of the pRF model as the compressive spatial summation  
611 (CSS) model, following Kay et al. (2013b). Briefly, we estimated the receptive field parameters that, when applied to  
612 the drifting bar stimulus images, minimized the difference between the observed and predicted BOLD time series.  
613 First, stimulus images were converted to contrast apertures and downsampled to 101 x 101 grids. time series  
614 from each retinotopy run were resampled to anatomical space and restricted to gray matter voxels. time series  
615 were then averaged across runs. pRF models were solved using a two stage coarse-to-fine fit on the average  
616 time series. The first stage of the model fit was a coarse grid fit, which was used to find an approximate solution  
617 robust to local minima. This stage was solved on a volume-based time series that was first temporally decimated,  
618 spatially blurred on the cortical surface, and spatially subsampled. The parameters obtained with this fit were  
619 interpolated and then used as a seed for subsequent nonlinear optimization, or fine fit. This procedure yielded four  
620 final parameters of interest for every voxel: eccentricity ( $r$ ), polar angle ( $\theta$ ), sigma ( $\sigma$ ), exponent ( $n$ ). The eccentricity  
621 and polar angle parameters describe the location of the receptive field in space, the sigma parameter describes the  
622 size of the receptive field, and the exponent describes the amount of compressive spatial summation applied to  
623 responses from the receptive field. Eccentricity and polar angle parameters were converted from polar coordinates  
624 to rectangular coordinates ( $x, y$ ) for some analyses. Variance explained by the pRF model with these parameters  
625 was also calculated for each voxel. All parameters were then projected from each subject's anatomical volume to  
626 the cortical surface (Fig. 2b, bottom).

627 **ROI definitions**

628 Regions of interest were defined by hand-drawing boundaries at polar angle reversals on each subject's cortical  
629 surface, following established practice (Wandell et al., 2007). We used this method to define six ROIs spanning  
630 early to mid-level visual cortex: V1, V2, V3, hV4, LO (LO1 and LO2), and V3ab (V3a and V3b).

631 We further restricted each ROI by preferred eccentricity in order to isolate vertices responsive to our stimuli. We  
632 excluded vertices with eccentricity values less than 0.5° and greater than 8°. This procedure excluded vertices  
633 responding primarily to the fixation dot and vertices near the maximal extent of visual stimulation in the scanner. We  
634 also excluded vertices whose variance explained by the pRF model ( $R^2$ ) was less than 0.1, indicating poor spatial  
635 selectivity. All measures used to exclude vertices from ROIs were independent of the measurements made during  
636 the perception and memory tasks.

637 **Analyses quantifying perception and memory activity**

638 Our main empirical analyses examined the evoked BOLD response to our experimental stimuli during perception  
639 and memory as a function of visual field parameters estimated from the pRF model. Our first step was to visualize  
640 evoked activity during perception and memory in visual field coordinates (Fig. 3a). Transforming the data in this  
641 way allowed us to view the activity in a common reference frame across all brain regions, rather than on the  
642 cortical surface, where comparisons are made difficult by the fact that surface area and cortical magnification differ  
643 substantially from one area to the next. To do this, we selected the ( $x, y$ ) parameters for each surface vertex from the  
644 retinotopy model and the  $\beta$  parameters from the GLM analysis. Separately for a given ROI, subject, stimulus, and  
645 task (perception/memory), we interpolated the  $\beta$  values over ( $x, y$ ) space. We rotated each of these representations

according to the polar angle location of the stimulus so that they would be aligned at the upper vertical meridian. We then z-scored each representation before averaging across stimuli and subjects. We used these images to gain intuition about the response profiles and to guide subsequent quantitative analyses.

Before quantifying these representations, we simplified them further. Because our stimuli were all presented at the same eccentricity, we reduced our 2D stimulus coordinate representations to 1D dimensional responses functions on the polar angle dimension (Fig. 4a). We did this by selecting surface vertices whose  $(x, y)$  coordinates were within one  $\sigma$  of the stimulus eccentricity ( $2^\circ$ ) for each ROI. We then binned the evoked BOLD response into 18 bins of polar angle distance from the stimulus and averaged within each bin to produce polar angle response functions for each subject. We divided each subject's response function by the norm of the response vector before averaging across subjects and then multiplying by the average vector norm to get the correct units back. This procedure prevents a subject with a high BOLD response across all polar angles from dominating the average response. The resulting average polar angle response functions showed clear surround suppression for polar angles near the stimulus during perception. Given this, we fit a difference of two von Mises distributions to the average data, with the location parameters ( $\mu$ ) for the two von Mises distributions fixed to be equal, but the spread ( $\kappa$ ) and scale allowed to differ.

We quantitatively assessed the similarities and differences between perception and memory responses using these fits. We examined the location parameter of the two von Mises distributions, and also computed the amplitude and FWHM of the fit. We repeated the fitting procedure 500 times, drawing subjects with replacement, to create bootstrapped 68% and 95% confidence intervals for both perception and memory location, amplitude, and FWHM parameters. To assess main effects of ROI, main effects of perception vs memory, and the interaction of these variables on location, amplitude, and FWHM values, we ran two-way ANOVAs. In all models, ROI was coded as an ordinal variable ( $V1 < V2 < V3 < hV4 < LO < V3ab$ ) and perception/memory as a categorical variable. Because location, amplitude, and FWHM, were computed at the group-level and not at the single-subject level, we ran these ANOVAs using group-level values. We re-ran the ANOVAs for all 500 subject resamplings to create bootstrapped confidence intervals for ANOVA regression coefficients. We computed p-values for these effects by performing randomization tests. To create null distributions, we randomly shuffled the assignment of the location, amplitude, or FWHM values with respect to the independent variables of interest (ROI, perception/memory). We did this for every possible shuffling or a subset of 10,000 different shufflings, whichever was smaller. We then computed two-tailed p-values according to the position of the true regression coefficient in the null distribution. Statistical data visualizations for these analyses and those subsequently described were made using seaborn v0.9.0 (Waskom et al., 2018).

## Noise simulations

We performed three simulations designed to test whether differences in noise between perception and memory data could explain differences in the responses we observed. To this end, we identified three potential types of noise that were present in our memory data but not our perception data: 1) reduced SNR; 2) retrieval failure; 3) memory error. We then simulated the effect of these types of noise on our perception data and asked whether these noise sources could produce responses similar to the ones we observed during memory.

### **SNR simulation**

To simulate reduced SNR, we created artificial datasets with different amounts of additive noise introduced to every vertex's perception parameter estimate. Noise was added in five levels: noise needed to generate the empirical SNR of the perception data ( $p$ ), noise needed to generate the empirical SNR of the memory data ( $m$ ), or noise needed to generate 1/2, 1/4, or 1/8 the empirical SNR of the memory data. For each of these values, we simulated 100 independent datasets for every subject and ROI. We determined the amount of signal and noise actually observed for each vertex during perception and memory by examining bootstrapped parameter estimate distributions produced by GLMdenoise. We defined the median parameter estimate across bootstraps as the amount of signal and the standard error of this distribution as the amount of noise. To simulate new data for a vertex, we randomly drew a new parameter estimate from a normal distribution defined by the true signal value

693 (median) and the noise value (SE) needed to produce the target SNR. Critically, we made the draws correlated  
694 across vertices for each simulation. We did this by selecting a scale factor from a standard normal distribution  
695 which determined how many SEs away from the median every vertex's simulated value would lie. This scale factor  
696 was shared across all vertices in an ROI for a given simulation. This procedure overcompensates for the spatial  
697 correlation present in BOLD data by assuming that SNR is 100% correlated across all vertices in an ROI. Note  
698 that if the noise were uncorrelated across vertices, it would have a much smaller effect on the population tuning  
699 curves. For each noise value and each of the 100 simulations, we analyzed the simulated data using the same  
700 procedure we applied to the actual data. This yielded 100 von Mises fits to the simulated data for each noise value  
701 and ROI (Supplementary Fig. 1a). We extracted the location, amplitude, and FWHM values from these fits. We  
702 evaluated whether location and FWHM values approximated the ones we observed during memory by calculating  
703 the proportion of simulations that fell within the 95% confidence intervals derived from the memory data (Fig. 5a).

#### 704 **Retrieval failure simulation**

705 To simulate retrieval failure, we created artificial datasets that contained a variable number of perception trials  
706 with no signal. Retrieval failure was simulated in five levels: 0%, 25%, 50%, 75%, and 100% of trials. For each  
707 of these values, we simulated 100 independent datasets for every subject and ROI. Depending on the level of  
708 retrieval failure, zero, one, two, three, or all four stimuli were randomly designated as 'failed' in each simulated  
709 dataset. For the failed stimuli, new parameter estimates were drawn from a distribution defined by zero signal during  
710 perception for every vertex. For the remaining stimuli, new parameter estimates were drawn from a distribution  
711 defined by the true perception signal for every vertex. Noise was equated for both trial types; for each vertex, we  
712 used the the amount of noise observed during perception. As in the SNR simulation, simulated data were correlated  
713 across vertices in an ROI and simulated data were analyzed using the same procedures as for the actual data. We  
714 evaluated whether simulated location and FWHM values approximated the ones we observed during memory by  
715 calculating the proportion of simulations that fell within the 95% confidence intervals derived from the memory data  
(Supplementary Fig. 1b and Fig. 5b).

#### 717 **Memory error simulation**

718 To simulate memory error, we created artificial datasets that contained a variable amount of angular error in the peak  
719 location of the perception polar angle response functions. Memory error was simulated in seven levels of standard  
720 deviation: 0, 15, 30, 45, 60, 75, and 90 degrees. For each of these values, we simulated 100 independent datasets  
721 for every subject and ROI. We assigned the amount of memory error for a given subject and stimulus by drawing a  
722 random value from a normal distribution centered at the true angular location of the stimulus and with the current  
723 standard deviation. We then used these memory error values to misalign simulated perception data. Specifically,  
724 we created new perception datasets based on the *true* signal and noise characteristics of our perception data  
725 (equivalent to SNR simulation with 'p' noise or 0% retrieval failure simulation). As in prior simulations, simulated data  
726 were correlated across vertices in an ROI, and simulated data were analyzed according to the same procedure as  
727 for the actual data. Before averaging the simulated data across stimuli and subjects, we rotated each response by  
728 the chosen memory error value rather than by the location of that stimulus. That is, instead of rotating the response  
729 to a 45° stimulus by 45° to align all stimuli at 0° (as we did in our main analysis), we rotated the response by a  
730 value either close (generating using small standard deviations, representing small errors) or potentially quite far  
731 away (generating using large standard deviations, representing large errors). After averaging, we extracted location  
732 and FWHM values. We then evaluated whether simulated location and FWHM values approximated the ones we  
733 observed during memory by calculating the proportion of simulations that fell within the 95% confidence intervals  
734 derived from the memory data (Supplementary Fig. 1c and Fig. 5c).

#### 735 **pRF forward model**

736 We evaluated the ability of our pRF model to account for our perception and memory measurements. To do this, we  
737 used our pRF model as a forward model. This means that we took the pRF model parameters fit to fMRI data from  
738 the retinotopy session (which used a drifting bar stimulus) and used them to generate predicted BOLD responses  
739 to our four experimental stimuli. The model takes processed stimulus images as input, and for each of these

740 images, outputs a predicted BOLD response (in units of % signal change) for every cortical surface vertex. Before  
741 running the model, we transformed our experimental stimuli into binary contrast apertures with values of 1 where  
742 the stimulus was and values of 0 everywhere else. These images were downsampled to the same resolution as the  
743 images used to fit the pRF model (101 x 101).

744 **Model specification**

745 The pRF forward model has two fundamental operations. In the first operation, a stimulus contrast aperture image is  
746 multiplied by a voxel's pRF. In the CSS and linear models, this pRF is defined as a circular symmetric 2D Gaussian,  
747 parameterized by a location in the visual field ( $x, y$ ) and a size ( $\sigma$ ). In the DoG+CSS version of the model, this pRF  
748 is defined as the difference of two such Gaussians, centered at the same location (see next paragraph). The second  
749 operation applies a power-law exponent ( $n$ ) to the result of the multiplication, effectively boosting small responses.  
750 This nonlinear operation is the key component of the CSS model and improves model accuracy in high-level visual  
751 areas that are known to exhibit subadditive spatial summation (Kay et al., 2013b; Mackey et al., 2017). The values  
752 of the exponent range from 0 to 1, where a value of 1 returns the model to linear. The output of this nonlinear stage  
753 is multiplied by a final scale parameter ( $\beta$ ), which returns the units to % signal change (Fig. 6a).

754 Because we observed negative surround responses in V1–V3 during perception, we focused mainly on the  
755 results of the DoG+CSS model. Prior work has shown that difference-of-Gaussians (DoG) pRF models can account  
756 for the center-surround structure we observed (Zuiderbaan et al., 2012). In order to construct DoG pRFs, we  
757 converted each pRF from the CSS model we fit to the retinotopy data to a DoG pRF. We chose this approach after  
758 encountering difficulty in fitting a DoG pRF model to the retinotopy data. First, we took every 2D Gaussian pRF from  
759 the CSS model, and we subtracted from it a second 2D Gaussian pRF that was centered at the same location but  
760 was twice as wide and half as high. This ratio of  $2\sigma$  and  $.5\beta$  between the negative and positive Gaussians was fixed  
761 for all voxels. In order to prevent the resulting DoG pRF from being systematically narrower and lower in amplitude  
762 than the original pRF, we rescaled the  $\sigma$  and  $\beta$  of the original pRF before converting it to a DoG. We multiplied  
763 the original  $\sigma$  by  $\sqrt{2}$  and the original  $\beta$  by 2, resulting in a DoG pRF with equivalent FWHM and amplitude as the  
764 original pRF. Thus, the DoG pRF differed from the original pRF only in the presence of a suppressive surround.

765 We compared the predicts of the DoG+CSS model to the results of the CSS model and to a linear model that we  
766 fit separately to the retinotopy data. In this linear model, no exponent parameter was fit. After generating a prediction  
767 for each subject, stimulus, and surface vertex, for each of our three forward models, we carried these predictions  
768 forward through the same analysis pipeline used to analyze our task-based data. This generated predicted polar  
769 angle response functions for each of the three pRF forward models (Fig. 6b and Supplementary Fig. 2). We  
770 generated bootstrapped predictions by conducting the same procedure on the bootstrapped datasets.

771 **Evaluating model predictions**

772 We next compared how well the DOG+CSS model predictions matched our perception versus memory measure-  
773 ments. We correlated the predicted location, amplitude, and FWHM parameters for each ROI with the actual  
774 perception and memory parameters. We evaluated these relationships by fitting a linear model to the predicted  
775 versus observed observations. To generate confidence intervals on these fits, we fit linear models to the 500  
776 bootstrapped perception and memory datasets and the yoked pRF predictions (Fig. 6c).

777 We also compared the model accuracy of the DoG+CSS and CSS predictions alongside a linear prediction  
778 with no exponent parameter (Supplementary Fig. 2a). We calculated the coefficient of determination ( $R^2$ ) for the  
779 predicted polar angle response functions in each ROI, separately for the observed perception and memory polar  
780 angle response functions (Supplementary Fig. 2b). Under this measure, a model that predicts the mean observed  
781 response for every value of polar angle distance will have an  $R^2$  of zero, with better models producing positive values  
782 and worse models producing negative values. We generated confidence intervals for these accuracies by computing  
783  $R^2$  values for each of the 500 bootstrapped perception and memory datasets and the yoked pRF predictions.

784 **Hierarchical network model**

785 We assessed whether a simple instantiation of a single neural network model could account for both the perception  
786 and memory data. We implemented a fully linear hierarchical model of neocortex in which the activity from each  
787 layer was created by pooling activity from the previous layer. This model encodes 1D space only and its parameters  
788 are fixed (i.e. it is not trained). For the feedforward simulation, we began with a 1D square wave stimulus, which  
789 spanned -20 to 20 degrees of polar angle. We created a fixed Gaussian convolution kernel ( $\mu = 0$ ,  $\sigma = 15$ ), which  
790 we convolved with the stimulus to create the activity in layer 1. This layer 1 activity was convolved with the same  
791 Gaussian kernel to create the layer 2 activity, and this process was repeated recursively for 8 layers (Fig. 7a, left).  
792 In order to simulate memory-evoked responses in this network, we made two assumptions. First, we assumed  
793 that the feedback simulation began with the layer 8 activity from the feedforward simulation. That is, we assumed  
794 no information loss or distortion between perception and memory in the last layer. Second, we assumed that  
795 all connections were reciprocal and thus that the same Gaussian kernel was applied to transform layers in the  
796 feedback direction as in the feedforward direction (Fig. 7a, right). Thus, in the feedback simulation, we convolved  
797 the layer 8 activity with the Gaussian kernel to produce the layer 7 activity and repeated this procedure recursively,  
798 ending at layer 1 (Fig. 7b). Note that these computations can be performed with matrix multiplication rather than  
799 convolution by converting the convolutional kernel to a Toeplitz matrix, which is how we implemented it. In this case,  
800 the transpose of the Toeplitz matrix (itself, as it is symmetric) is used in the feedback direction. We plot the location,  
801 amplitude and FWHM for each layer's activation in the same convention as the data (Fig. 7c).

802 **Acknowledgements**

803 S.E.F. was supported by an NIH Blueprint D-SPAN F99/K00 Award (F99-NS105223).

804 **Author Contributions**

805 Conceptualization, S.E.F., B.A.K., and J.W.; Methodology, S.E.F. and J.W.; Software, S.E.F. and J.W.; Investigation,  
806 S.E.F.; Writing–Original Draft, S.E.F. and J.W.; Writing–Review & Editing, S.E.F., B.A.K., and J.W.; Funding  
807 Acquisition, S.E.F.; Supervision, B.A.K and J.W.

808 **Declaration of Interests**

809 The authors declare no competing interests.

810 **References**

- 811 Barone, P., Batardiere, A., Knoblauch, K., & Kennedy, H. (2000). Laminar distribution of neurons in extrastriate  
812 areas projecting to visual areas V1 and V4 correlates with the hierarchical rank and intimates the operation of a  
813 distance rule. *Journal of Neuroscience*, 20(9), 3263–3281.
- 814 Benson, N. C., Jamison, K. W., Arcaro, M. J., Vu, A. T., Glasser, M. F., Coalson, T. S., Van Essen, D. C., Yacoub, E.,  
815 Ugurbil, K., Winawer, J., & Kay, K. (2018). The Human Connectome Project 7 Tesla retinotopy dataset: Description  
816 and population receptive field analysis. *Journal of Vision*, 18(13), 23.
- 817 Benson, N. C. & Winawer, J. (2018). Bayesian analysis of retinotopic maps. *eLife*, 7, 0–45.
- 818 Bloem, I. M., Watanabe, Y. L., Kibbe, M. M., & Ling, S. (2018). Visual Memories Bypass Normalization. *Psychological  
819 science*, 29(5), 845–856.
- 820 Bone, M. B., St-Laurent, M., Dang, C., McQuiggan, D. A., Ryan, J. D., & Buchsbaum, B. R. (2018). Eye Movement  
821 Reinstatement and Neural Reactivation During Mental Imagery. *Cerebral Cortex*, 29(3), 1075–1089.

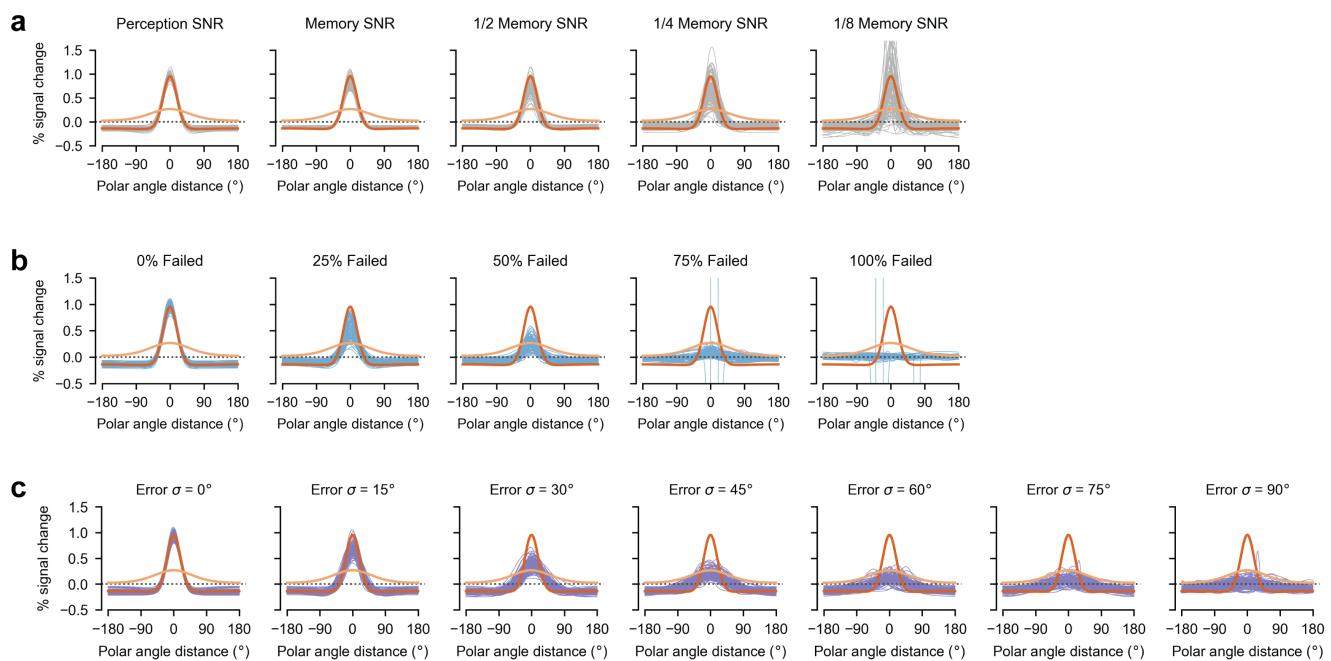
- 822 Bosch, S. E., Jehee, J. F. M., Fernandez, G., & Doeller, C. F. (2014). Reinstatement of Associative Memories in  
823 Early Visual Cortex Is Signaled by the Hippocampus. *Journal of Neuroscience*, 34(22), 7493–7500.
- 824 Breedlove, J. L., St-Yves, G., Olman, C. A., & Naselaris, T. P. (2018). Mental imagery encoding models reveal  
825 signatures of inference in a hierarchical generative model. *bioRxiv*.
- 826 Buracas, G. T. & Boynton, G. M. (2007). The Effect of Spatial Attention on Contrast Response Functions in Human  
827 Visual Cortex. *Journal of Neuroscience*, 27(1), 93–97.
- 828 Carrasco, M. (2011). Visual attention: The past 25 years. *Vision Research*, 51(13), 1484–1525.
- 829 Chawla, D., Rees, G., & Friston, K. J. (1999). The physiological basis of attentional modulation in extrastriate visual  
830 areas. *Nature Neuroscience*, 2(7), 671–676.
- 831 Damasio, A. R. (1989). Time-locked multiregional retroactivation: A systems level proposal for the neural substrates  
832 of recall and recognition. *Cognition*, 33, 25–62.
- 833 Dijkstra, N., Ambrogioni, L., & Gerven, M. A. J. V. (2019). Neural dynamics of perceptual inference and its reversal  
834 during imagery. *bioRxiv*.
- 835 Dougherty, R. F., Koch, V. M., Brewer, A. A., Fischer, B., Modersitzki, J., & Wandell, B. A. (2003). Visual field  
836 representations and locations of visual areas v1/2/3 in human visual cortex. *Journal of Vision*, 3(10), 586–598.
- 837 Dumoulin, S. O. & Wandell, B. A. (2008). Population receptive field estimates in human visual cortex. *NeuroImage*,  
838 39(2), 647–660.
- 839 Engel, S. A., Rumelhart, D. E., Wandell, B. A., Lee, A. T., Glover, G. H., Chichilnisky, E.-J., & Shadlen, M. N. (1994).  
840 fMRI of human visual cortex.
- 841 Ester, E. F., Anderson, D. E., Serences, J. T., & Awh, E. (2013). A Neural Measure of Precision in Visual Working  
842 Memory. *Journal of Cognitive Neuroscience*, 25(5), 754–761.
- 843 Favila, S. E., Samide, R., Sweigart, S. C., & Kuhl, B. A. (2018). Parietal representations of stimulus features  
844 are amplified during memory retrieval and flexibly aligned with top-down goals. *Journal of Neuroscience*, 38(36),  
845 0564–18.
- 846 Felleman, D. J. & Essen, D. C. V. (1991). Distributed Hierarchical Processing in the Primate Cerebral Cortex.  
847 *Cerebral Cortex*, 1, 1–47.
- 848 Fischl, B. (2012). FreeSurfer. *NeuroImage*, 62(2), 774–781.
- 849 Fukushima, K. (1980). Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition  
850 unaffected by shift in position. *Biological Cybernetics*, 36(4), 193–202.
- 851 Gandhi, S. P., Heeger, D. J., & Boynton, G. M. (1999). Spatial attention affects brain activity in human primary. *Proc.  
852 Natl. Acad. Sci. USA*, 96(March), 3314–3319.
- 853 Gattass, R., Nascimento-Silva, S., Soares, J. G., Lima, B., Jansen, A. K., Diogo, A. C. M., Farias, M. F., Botelho,  
854 Eliā P, M. M., Mariani, O. S., Azzi, J., & Fiorani, M. (2005). Cortical visual areas in monkeys: location, topography,  
855 connections, columns, plasticity and cortical dynamics. *Philosophical Transactions of the Royal Society B: Biological  
856 Sciences*, 360(1456), 709–731.
- 857 Gordon, A. M., Rissman, J., Kiani, R., & Wagner, A. D. (2014). Cortical Reinstatement Mediates the Relationship  
858 Between Content-Specific Encoding Activity and Subsequent Recollection Decisions. *Cerebral Cortex*, 24(12),  
859 3350–3364.
- 860 Gorgolewski, K., Madison, C., Burns, C. D., Clark, D., Halchenko, Y. O., Waskom, M. L., & Ghosh, S. S. (2011).  
861 NiPy: A Flexible, Lightweight and Extensible Neuroimaging Data Processing Framework in Python. *Frontiers in  
862 Neuroinformatics*, 5(August).

- 863 Hebb, D. O. (1968). Concerning imagery. *Psychological Review*, 75(6), 466–77.
- 864 Heeger, D. J. (2017). Theory of cortical function. *Proceedings of the National Academy of Sciences*, 114(8),  
865 1773–1782.
- 866 Hilgetag, C. C. & Goulas, A. (2020). ‘Hierarchy’ in the organization of brain networks. *Philosophical Transactions of  
867 the Royal Society B: Biological Sciences*, 375(1796), 20190319. Publisher: Royal Society.
- 868 Hindy, N. C., Ng, F. Y., & Turk-Browne, N. B. (2016). Linking pattern completion in the hippocampus to predictive  
869 coding in visual cortex. *Nature Neuroscience*, 19(5), 665–667. Publisher: Nature Publishing Group ISBN: 1546-1726.
- 870 Horikawa, T. & Kamitani, Y. (2017). Generic decoding of seen and imagined objects using hierarchical visual  
871 features. *Nature Communications*, 8(May), 1–15.
- 872 Ishai, A. & Sagi, D. (1995). Common mechanisms of visual imagery and perception. *Science*, 268(5218), 1772–1774.
- 873 James, W. (1890). *The Principles of Psychology*. New York, NY: Holt.
- 874 Kastner, S., Pinsk, M. A., De Weerd, P., Desimone, R., & Ungerleider, L. G. (1999). Increased activity in human  
875 visual cortex during directed attention in the absence of visual stimulation. *Neuron*, 22(4), 751–61.
- 876 Kay, K. N., Rokem, A., Winawer, J., Dougherty, R. F., & Wandell, B. A. (2013a). GLMdenoise: A fast, automated  
877 technique for denoising task-based fMRI data. *Frontiers in Neuroscience*, 7(7 DEC), 1–15.
- 878 Kay, K. N., Winawer, J., Mezer, A., & Wandell, B. A. (2013b). Compressive spatial summation in human visual  
879 cortex. *Journal of Neurophysiology*, 110(2), 481–494.
- 880 Kosslyn, S. M., Thompson, W. L., Kim, I. J., & Alpert, N. M. (1995). Topographical representations of mental images  
881 in primary visual cortex. *Nature*, 378(6556), 496–8.
- 882 Kuhl, B. A., Johnson, M. K., & Chun, M. M. (2013). Dissociable neural mechanisms for goal-directed versus  
883 incidental memory reactivation. *The Journal of Neuroscience*, 33(41), 16099–109.
- 884 Kuhl, B. A., Rissman, J., Chun, M. M., & Wagner, A. D. (2011). Fidelity of neural reactivation reveals competition  
885 between memories. *Proceedings of the National Academy of Sciences*, 108(14), 5903–5908.
- 886 Lee, S. H., Kravitz, D. J., & Baker, C. I. (2012). Disentangling visual imagery and perception of real-world objects.  
887 *NeuroImage*, 59(4), 4064–4073.
- 888 Lee, S.-h., Kravitz, D. J., & Baker, C. I. (2018). Differential Representations of Perceived and Retrieved Visual  
889 Information in Hippocampus and Cortex. *Cerebral Cortex*, (pp. 1–10).
- 890 Li, X., Lu, Z.-L., Tjan, B. S., Dosher, B. A., & Chu, W. (2008). Blood oxygenation level-dependent contrast response  
891 functions identify mechanisms of covert attention in early visual areas. *Proceedings of the National Academy of  
892 Sciences*, 105(16), 6202–6207.
- 893 Linde-Domingo, J., Treder, M. S., Kerrén, C., & Wimber, M. (2019). Evidence that neural information flow is reversed  
894 between object perception and object reconstruction from memory. *Nature Communications*, 10(1), 179.
- 895 Luck, S. J., Chelazzi, L., Hillyard, S. A., & Desimone, R. (1997). Neural Mechanisms of Spatial Selective Attention  
896 in Areas V1, V2, and V4 of Macaque Visual Cortex. *Journal of Neurophysiology*, 77(1), 24–42.
- 897 Mackey, W. E., Winawer, J., & Curtis, C. E. (2017). Visual field map clusters in human frontoparietal cortex. *eLife*,  
898 6(e22974).
- 899 Marr, D. (1971). Simple Memory: A Theory for Archicortex. *Philosophical Transactions of the Royal Society B:  
900 Biological Sciences*, 262(841), 23–81.
- 901 McClelland, J. L., McNaughton, B. L., & O'Reilly, R. C. (1995). Why there are complementary learning systems in  
902 the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and  
903 memory. *Psychological review*, 102(3), 419–457.

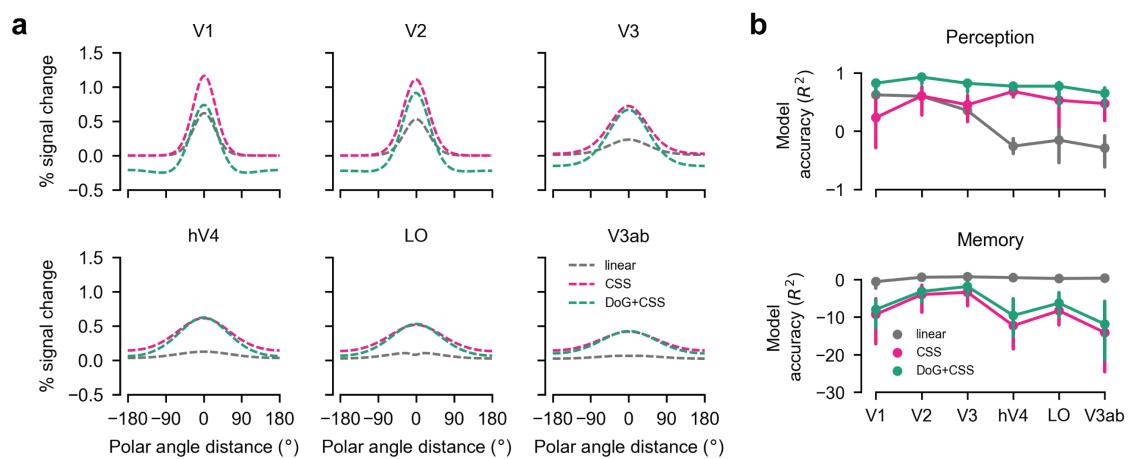
- 904 Naselaris, T., Kay, K. N., Nishimoto, S., & Gallant, J. L. (2011). Encoding and decoding in fMRI. *NeuroImage*, 56(2),  
905 400–410.
- 906 Naselaris, T., Olman, C. A., Stansbury, D. E., Ugurbil, K., & Gallant, J. L. (2015). A voxel-wise encoding model for  
907 early visual areas decodes mental images of remembered scenes. *NeuroImage*, 105, 215–228.
- 908 Naya, Y., Yoshida, M., & Miyashita, Y. (2001). Backward Spreading of Memory-Retrieval Signal in the Primate  
909 Temporal Cortex. *Science*, 291(5504), 661–664.
- 910 O’Craven, K. M. & Kanwisher, N. (2000). Mental Imagery of Faces and Places Activates Corresponding Stimulus-  
911 Specific Brain Regions. *Journal of Cognitive Neuroscience*, 12(6), 1013–1023.
- 912 O'Reilly, R. C. & McClelland, J. L. (1994). Hippocampal conjunctive encoding, storage, and recall: Avoiding a  
913 trade-off. *Hippocampus*, 4(6), 661–682.
- 914 Pearson, J. (2019). The human imagination: the cognitive neuroscience of visual mental imagery. *Nature Reviews  
915 Neuroscience*.
- 916 Pearson, J., Clifford, C. W., & Tong, F. (2008). The Functional Impact of Mental Imagery on Conscious Perception.  
917 *Current Biology*, 18(13), 982–986.
- 918 Pearson, J., Naselaris, T., Holmes, E. A., & Kosslyn, S. M. (2015). Mental Imagery: Functional Mechanisms and  
919 Clinical Applications. *Trends in Cognitive Sciences*, 19(10), 590–602.
- 920 Polyn, S. M., Natu, V. S., Cohen, J. D., & Norman, K. A. (2005). Category-Specific Cortical Activity Precedes  
921 Retrieval During Memory Search. *Science*, 310(5756), 1963–6.
- 922 Pylyshyn, Z. W. (2002). Mental imagery: In search of a theory. *Behavioral and Brain Sciences*, 25(2), 157–182.
- 923 Rademaker, R. L., Chunharas, C., & Serences, J. T. (2019). Coexisting representations of sensory and mnemonic  
924 information in human visual cortex. *Nature Neuroscience*, 22(8). Publisher: Springer US.
- 925 Rahmati, M., Saber, G., & Curtis, C. (2017). Population Dynamics of Early Visual Cortex During Working Memory.  
926 *Journal of Cognitive Neuroscience*.
- 927 Ress, D., Backus, B. T., & Heeger, D. J. (2000). Activity in primary visual cortex predicts performance in a visual  
928 detection task. *Nature Neuroscience*, 3(9), 940–945.
- 929 Riesenhuber, M. & Poggio, T. (1999). Hierarchical models of object recognition in cortex. *Nature Neuroscience*,  
930 2(11), 1019–1025.
- 931 Sereno, M. I., Dale, A. M., Reppas, J. B., Kwong, K. K., Belliveau, J. W., Brady, T. J., Rosen, B. R., & Tootell, R. B.  
932 (1995). Borders of multiple visual areas in humans revealed by functional magnetic resonance imaging. *Science*,  
933 268(5212), 889–893.
- 934 Serre, T., Oliva, A., & Poggio, T. (2007). A feedforward architecture accounts for rapid categorization. *Proceedings  
935 of the National Academy of Sciences*, 104(15), 6424–6429.
- 936 Slotnick, S. D., Thompson, W. L., & Kosslyn, S. M. (2005). Visual mental imagery induces retinotopically organized  
937 activation of early visual areas. *Cerebral Cortex*, 15(10), 1570–1583.
- 938 Smith, S. M., Jenkinson, M., Woolrich, M. W., Beckmann, C. F., Behrens, T. E., Johansen-Berg, H., Bannister, P. R.,  
939 Luca, M. D., Drobniak, I., Flitney, D. E., Niazy, R. K., Saunders, J., Vickers, J., Zhang, Y., Stefano, N. D., Brady, J. M.,  
940 & Matthews, P. M. (2004). Advances in functional and structural MR image analysis and implementation as FSL.  
941 *NeuroImage*, 23, S208–S219.
- 942 Somers, D. C., Dale, A. M., Seiffert, A. E., & Tootell, R. B. H. (1999). Functional MRI reveals spatially specific  
943 attentional modulation in human primary visual cortex. *Proceedings of the National Academy of Sciences*, 96(4),  
944 1663–1668.

- 945 Sperling, G. (1960). The information available in brief visual presentations. *Psychological Monographs: General  
946 and Applied*, 74(11), 1–29.
- 947 Sprague, T. C., Ester, E. F., & Serences, J. T. (2014). Reconstructions of information in visual spatial working  
948 memory degrade with memory load. *Current Biology*, 24(18), 2174–2180.
- 949 Sprague, T. C. & Serences, J. T. (2013). Attention modulates spatial priority maps in the human occipital, parietal  
950 and frontal cortices. *Nature neuroscience*, 16(12), 1879–1887.
- 951 Sutterer, D. W., Foster, J. J., Serences, J. T., Vogel, E. K., & Awh, E. (2019). Alpha-band oscillations track the  
952 retrieval of precise spatial representations from long-term memory. *Journal of Neurophysiology*, 122(2), 539–551.
- 953 Suzuki, W. A. & Amaral, D. G. (1994). Perirhinal and parahippocampal cortices of the macaque monkey: Cortical  
954 afferents 4025. *Journal of Comparative Neurology*, 350, 497–533.
- 955 Tartaglia, E. M., Bamert, L., Mast, F. W., & Herzog, M. H. (2009). Human Perceptual Learning by Mental Imagery.  
956 *Current Biology*, 19(24), 2081–2085.
- 957 Thirion, B., Duchesnay, E., Hubbard, E., Dubois, J., Poline, J. B., Lebihan, D., & Dehaene, S. (2006). Inverse  
958 retinotopy: Inferring the visual content of images from brain activation patterns. *NeuroImage*, 33(4), 1104–1116.
- 959 Van Hoesen, G. & Pandya, D. N. (1975). Some connections of the entorhinal (area 28) and perirhinal (area 35)  
960 cortices of the rhesus monkey. I. Temporal lobe afferents. *Brain Research*, 95(1), 1–24.
- 961 Waldhauser, G. T., Braun, V., & Hanslmayr, S. (2016). Episodic Memory Retrieval Functionally Relies on Very Rapid  
962 Reactivation of Sensory Information. *The Journal of Neuroscience*, 36(1), 251–260.
- 963 Wandell, B., Dumoulin, S. O. S., & Brewer, A. A. a. (2007). Visual Field Maps in Human Cortex. *Neuron*, 56(2),  
964 366–383.
- 965 Wandell, B. A. & Winawer, J. (2015). Computational neuroimaging and population receptive fields. *Trends in  
966 Cognitive Sciences*, 19(6), 349–357.
- 967 Waskom, M., Botvinnik, O., O’Kane, D., Hobson, P., Ostblom, J., Lukauskas, S., Gemperline, D. C., Augspurger, T.,  
968 Halchenko, Y., Cole, J. B., Warmenhoven, J., de Ruiter, J., Pye, C., Hoyer, S., Vanderplas, J., Villalba, S., Kunter, G.,  
969 Quintero, E., Bachant, P., Martin, M., Meyer, K., Miles, A., Ram, Y., Brunner, T., Yarkoni, T., Williams, M. L., Evans,  
970 C., Fitzgerald, C., Brian, & Qalieh, A. (2018). mwaskom/seaborn: v0.9.0 (july 2018).
- 971 Wheeler, M. E., Petersen, S. E., & Buckner, R. L. (2000). Memory’s echo: Vivid remembering reactivates sensory-  
972 specific cortex. *Proceedings of the National Academy of Sciences*, 97(20), 11125–11129.
- 973 Winawer, J., Huk, A. C., & Boroditsky, L. (2010). A motion aftereffect from visual imagery of motion. *Cognition*,  
974 114(2), 276–284.
- 975 Zeki, S. (2015). A massively asynchronous, parallel brain. *Philosophical Transactions of the Royal Society B:  
976 Biological Sciences*, 370(1668), 20140174. Publisher: Royal Society.
- 977 Zuiderbaan, W., Harvey, B. M., & Dumoulin, S. O. (2012). Modeling center-surround configurations in population  
978 receptive fields using fMRI. *Journal of Vision*, 12(3), 10–10.

979 **Supplementary Figures**



**Figure 1. Simulated V1 datasets with different noise levels.** (a) Gray lines represent the fits to simulated V1 perception datasets with different levels of SNR. Each panel contains 100 independently simulated datasets with the same noise level. Orange lines represent the fits to the actual perception and memory data, reproduced from Figure 4b, and are the same for each SNR value. (b) Purple lines represent the fits to simulated V1 perception datasets with different frequencies of failed retrieval. Other conventions as in (a). (c) Blues lines represent the fits to simulated V1 perception datasets with different amounts of memory error. Other conventions as in (a).



**Figure 2. pRF model comparisons.** (a) Predicted polar angle response functions are plotted for three pRF models: linear, CSS, and DoG+CSS. Comparing these responses to perception data plotted in Figure 4b, the linear model did the poorest job of predicting perception responses. Linear predictions underestimated the amplitude of the observed response, particularly in later visual areas. Both nonlinear models (CSS and DOG+CSS) avoided this magnitude of failure. The DoG+CSS model selectively captured negative responses in V1–V3. (b) Model accuracy ( $R^2$ ) of the predicted polar angle response functions for each pRF model, evaluated separately for perception and memory data in each ROI. Error bars indicate 68% bootstrapped confidence intervals. Accuracy of the linear model in predicting perception data dropped steadily moving away from V1, indicating poor fit. Model accuracies for the the CSS and DoG+CSS models were higher and more stable across ROIs, with the DoG+CSS performing slightly better in every region. With the exception of the linear model in late visual areas, accuracy for all three models was far worse for memory data than perception data.