1  **Title: The neural codes underlying internally generated representations in visual working**

2  **memory**

3  **Short Title: Internally generated neural codes in memory**

4

5  **Authors**:

6  Qing Yu[1], Bradley R. Postle[2,3]

7

8  **Affiliations**:

9  [1]Center for Excellence in Brain Science and Intelligence Technology, Chinese Academy of

10  Sciences, Shanghai, China

11  [2]Department of Psychiatry, University of Wisconsin–Madison, Madison, WI 53719, USA

12  [3]Department of Psychology, University of Wisconsin–Madison, Madison, WI 53706, USA

13

14  **Corresponding author:**

15  Qing Yu

16  Email: qingyu@ion.ac.cn

17

23

24 **Abstract**

25   Humans can construct rich subjective experience even when no information is available

26 in the external world. Here we investigated the neural representation of purely internally

27 generated stimulus-like information during visual working memory. Participants performed

28 delayed recall of oriented gratings embedded in noise with varying contrast during fMRI

29 scanning. Their trialwise behavioral responses provided an estimate of their mental

30 representation of the to-be-reported orientation. We used multivariate inverted encoding models

31 to reconstruct the neural representations of orientation in reference to the response. We found

32 that response orientation could be successfully reconstructed from activity in early visual cortex,

33 even on 0% contrast trials when no orientation information was actually presented, suggesting

34 the existence of a purely internally-generated neural code in early visual cortex. Additionally,

35 cross-generalization and multidimensional scaling analyses demonstrated that information

36 derived from internal sources was represented differently from typical working memory

37 representations, which receive influences from both external and internal sources. Similar results

38 were also observed in intraparietal sulcus (IPS), with slightly different cross-generalization

39 patterns. These results suggest a potential mechanism for how externally-driven and internally-

40 generated information is maintained in working memory.

41

45

46

47

## Introduction

Humans have the ability to mentally retain and manipulate visual information even when the information is not in view. This ability -- visual working memory -- is fundamental to human cognition (Baddeley, 2003; Engle, Tuholski, Laughlin, & Conway, 1999; Luck & Vogel, 2013). Understanding how the brain keeps such information online is thus a critical question for cognitive neuroscience. The sensorimotor recruitment hypothesis posits that sensory cortex is an important substrate for the representation of fine-grained perceptual information in working memory (Awh & Jonides, 2001; D'Esposito & Postle, 2015; Serences, 2016); for example, early visual cortex for maintaining low-level visual information. This view is supported by evidence from multivariate analyses of functional magnetic resonance imaging (fMRI) data that stimulus-specific information can be decoded from early visual cortex during maintenance of visual feature information (Harrison & Tong, 2009; Riggall & Postle, 2012; Serences, Ester, Vogel, & Awh, 2009; Yu & Shim, 2017). With fMRI, a neural code is assessed as a systematic set of mappings between different values of a cognitive state and different patterns of fMRI activity, and a shared code is inferred if the same mapping is observed across two domains of cognition. With this logic it has been demonstrated that, in early visual cortex, visual working memory shares the same neural codes with visual perception (Harrison & Tong, 2009), attention (Yu & Shim, 2019), and imagery (Albers, Kok, Toni, Dijkerman, & de Lange, 2013), suggesting that early visual cortex may serve as a mental buffer for representing visual information across different categories of cognitive task (Roelfsema & de Lange, 2016).

Although early visual cortex recruits common neural codes for different cognitive processes, these processes can be driven by distinct sources of information. For example, visual perception is driven by external, bottom-up input received from the retina, and visual imagery is driven by internal, top-down input from higher cortical areas (Pearson, 2019). Of course, early

3

72    visual cortex contains numerous reciprocal connections with higher cortical areas, and bottom-up

73    and top-down signaling are involved in most, if not all visually mediated behaviors, including

74    visual perception (Gilbert & Li, 2013; Muckli & Petro, 2013). Nevertheless, the fact that visual

75    imagery shows distinct temporal dynamics, and evolves later in time, compared with visual

76    perception (Dijkstra, Mostert, Lange, Bosch, & van Gerven, 2018), suggests at least some

77    meaningful distinction between the processing of externally presented and internally generated

78    information.

79          When considering the sensorimotor recruitment hypothesis, it is important to note that

80    visual working memory cannot be understood as merely the prolongation of sensory processing,

81    because many stimulus-nonspecific factors can influence representations in working memory.

82    For example, several studies have demonstrated recall biases toward discrete color centers in

83    visual working memory for color (Bae, Olkkonen, Allred, & Flombaum, 2015; Bae, Olkkonen,

84    Allred, Wilson, & Flombaum, 2014; Panichello, DePasquale, Pillow, & Buschman, 2019),

85    probably due to drift towards stable attractor states established through prior experience

86    (Panichello et al., 2019). Information from the previous trial can also be reactivated or otherwise

87    influence the current trial (Bae & Luck, 2019; Barbosa et al., 2020). Moreover, there is

88    considerable physiological evidence for an important role for feedback from higher cortical

89    areas. For example, laminar recordings indicate that delay-period input to V1 is most prominent

90    in supra- and infragranular layers that receive feedback projections from higher areas (Lawrence

91    et al., 2018; van Kerkoerle, Self, & Roelfsema, 2017), and delay-period local field potentials in

92    area MT are coherent with spiking in prefrontal cortex (PFC; Mendoza-Halliday, Torres, &

93    Martinez-Trujillo, 2014). Because typical working memory tasks, including those cited here,

94    begin with the external presentation of to-be-remembered stimulus information, delay-period

95    representations presumably reflect the combined influence of processes associated with

96    externally presented and with internally generated information. How such "typical" working

97    memory representations may differ from purely internally generated representations is the focus

98    of the present study. Although our focus is on representations in early visual cortex, we also

99    present results from intraparietal sulcus (IPS), because this region has also been implicated in

100   representing working memory-related information (e.g., Bettencourt & Xu, 2016; Ester, Sprague,

101   & Serences, 2015; Yu & Shim, 2017).

102        In the current study, stimulus contrast varied across trials between 0%, 10%, and 60%,

103   but participants were instructed that a sample orientation would be presented on every trial, and

104   that a recall response was required at the end of every trial, regardless of subjective visibility.

105   This allowed us to use responses to infer backwards what they had represented during the delay

106   period, including on 0% contrast ("null") trials that lacked external input. These responses could

107   then be used to investigate internally generated representations in visual working memory

108   maintenance. The comparison between the null and typical working memory trials (with 10% or

109   60% contrast) could also be used to isolate processes specific to internally generated

110   representations.

111

112   **Materials and Methods**

113   **Participants**

114        All participants were recruited from the University of Wisconsin–Madison community.

115   Two behavioral experiments (Experiments 1A and 1B) were performed prior to the fMRI

116   experiment (Experiment 2) to test the visibility of the stimuli to be used in the fMRI experiment.

117   Thirteen individuals (2 males, mean age $21.0 \pm 3.3$ years) participated in Experiment 1A, and 7

118   of these also participated in Experiment 1B, along with 9 new individuals (n = 16 in total; 3

119   males, mean age $19.6 \pm 1.9$ years). Eighteen individuals (including one who also participated in

5

120    Experiment 1B) participated in Experiment 2. One of these was excluded due to failure to

121    comply with task instructions, resulting in 17 individuals (4 males, mean age $23.5 \pm 3.8$ years) as

122    the final sample size for Experiment 2. We did not carry out power analysis a priori, but our

123    sample size was comparable or even superior to those from recent fMRI studies that have used a

124    similar task design (Bettencourt & Xu, 2016; Ester et al., 2015; Rademaker, Chunharas, &

125    Serences, 2019; Yu, Teng, & Postle, 2020). All participants had normal or corrected-to-normal

126    vision, reported no neurological or psychiatric disease, and provided written informed consent

127    approved by the University of Wisconsin–Madison Health Sciences Institutional Review Board.

128    All were monetarily compensated for their participation.

129

130    **Stimuli and procedure**

131        Sample stimuli were sinusoidal gratings embedded in white noise (spatial frequency =

132    1º/cycle, radius = 4º), presented at varying levels of Michelson contrast. In Experiment 1A, there

133    were two types of stimuli: gratings with a *high* contrast (60%), and gratings with a 75%

134    threshold-level contrast, determined for each subject with a thresholding task. In Experiments 1B

135    and 2, there were three types of stimuli: gratings with a *high* contrast (contrast = 60%), gratings

136    with a *low* contrast (contrast = 10%), and *null* stimuli (contrast = 0%). Importantly, no

137    orientation information was visible in *null* gratings, making them equivalent to white noise

138    patches.

139        All stimuli were created and presented using MATLAB (MathWorks, Natick, MA) and

140    Psychtoolbox 3 extensions (Brainard, 1997; Pelli, 1997). In Experiments 1A and 1B, stimuli

141    were presented on a 21.5-inch iMac screen at a viewing distance of 63 cm and behavioral

142    responses were made with a computer mouse. In Experiment 2, stimuli were projected via a 60-

143    Hz projector (Avotec Silent Vision 6011, Avotec, Inc., Stuart, FL), and viewed through a coil-

144 mounted mirror in the MRI scanner at a viewing distance of 69 cm, and participants' behavioral

145 responses were made with an MR-compatible trackball response pad (Current Designs Inc.,

146 Philadelphia, PA). During the scan, eye position was monitored and recorded using the Avotec

147 RE-5700 eye-tracking system (Avotec, Inc., Stuart, FL).

148

149 *Experiment 2*

150       We begin with a detailed description of Experiment 2, the experiment of primary

151 theoretical interest, during which participants performed 1-item delayed recall of orientation in

152 the fMRI scanner. On each trial, participants viewed a sample stimulus (*high*, *low*, or *null*)

153 presented at the center of the screen for 0.5 s. After a delay of 9.5 s (or 8.5 s for two

154 participants), an orientation dial (radius = 4º) was presented centrally, and participants rotated the

155 dial until its needle matched the remembered orientation as precisely as possible in a 4-s

156 response window. Critically, participants were told that an oriented grating would be presented

157 on every trial, although its visibility would vary across trials, and they were instructed to make a

158 best guess when they were unsure about what the orientation was. Feedback (recall error) was

159 provided following the response period for 0.5 s, even on *null* trials, and recall error was

160 calculated as the angular difference between sample and response orientations, regardless of

161 whether or not the sample orientation had actually been visible (Figure 1). The sample

162 orientation for each trial was randomly selected from 1º to 180º in steps of 1º in the orientation

163 space. The starting position of the needle of the response dial was randomly chosen on every

164 trial, independent of the sample.

165

166                              <----- insert Figure 1 about here ----->

167

168    For four participants, total trial length was 22 s: for two of the four participants tested

169    with an 8.5-s delay (S01 and S02), the inter-trial interval (ITI) was 9 s, and for the two tested

170    with a 9.5-s delay (S03 and S04), ITI was 8 s. For all remaining participants, for whom the delay

171    was 9.5 s and ITI was 10 s, total trial length was 24 s. To match the number of time points across

172    participants, all analyses focused on the first 22 s of every trial.

173    Each run began with an 8-s fixation period, followed by 18 experimental trials, and the

174    ratio of trial types (*high*: *low*: *null*) during each run was 3:1:2 (i.e., 9 *high* trials, 3 *low* trials, and

175    6 *null* trials). For one participant (S01), the experimental run in the first scan session was

176    truncated to 12 trials (i.e., 6 *high* trials, 2 *low* trials, 4 *null* trials) due to a technical problem with

177    scanning. Each participant completed 28 to 32 runs across two scanning sessions. In total, twelve

178    participants completed 270 *high* trials, 90 *low* trials, and 180 *null* trials (S02, S05, S06, S08 to

179    S12, S14 to S17); two participants completed 288 *high* trials, 96 *low* trials, and 192 *null* trials

180    (S03 and S04); two completed 252 *high* trials, 84 *low* trials, and 168 *null* trials (S07 and S13);

181    and one (S01) completed 231 *high* trials, 77 *low* trials, and 154 *null* trials. All participants were

182    debriefed at the end of the study, and none of them reported awareness of the existence of null

183    trials (i.e., all reported believing that an oriented grating was presented on every trial).

184

185    *Experiments 1A and 1B*

186    Prior to the fMRI experiment, we ran two behavioral studies to determine the contrasts of

187    the gratings to be used in the scanner. The overarching rationale was to develop conditions that

188    would disguise from participants the fact that a substantial proportion of samples contained no

189    stimulus information (i.e., *null* samples). To achieve this, we sought to find two levels of contrast

190    that were each highly discriminable, but that would create the impression for participants that

191    subjective visibility would vary from trial to trial. The trial structure for both was similar to that

192    from Experiment 2: one sample grating (radius = 3º) with a randomly selected orientation was

193    presented on the screen for 0.1 s, followed by a brief delay, followed by recall with an

194    orientation wheel. Responses were self-paced, and feedback was given after each response (0.5

195    s).

196          Experiment 1A was carried out to examine how subjects would perform at each of two

197    levels of contrast: high and at-threshold. It began with a block of 80 trials to determine each

198    individual's contrast threshold: After an initial 10 trials of delayed recall (delay of 0.3 s) at a

199    fixed contrast of 12%, the sample contrast for each of the ensuing trials was adjusted using a

200    QUEST procedure (Watson & Pelli, 1983). Responses were binarized using a cut-off criterion of

201    20º of recall error. Four catch trials were interleaved at randomly determined intervals, and on

202    these catch trials the contrast was set to three times of the contrast from QUEST. The

203    discrimination contrast threshold of the grating that generated 75% accuracy was determined at

204    the end of the block. During the remainder of the session, participants performed 5 or 6 blocks of

205    delayed recall of orientation, delay length was either 1 s or 7 s, and delay length and sample

206    contrast (60%; at threshold) were fully crossed during each 60-trial block.

207          Experiment 1B was carried out to examine how subjects would perform at each of the

208    three levels of contrast that would be used for the fMRI study: *high* (60%); *low* (10%), and *null*

209    (0%). Participants performed 5 or 6 blocks with 60 trials each; again, delay length was either 1 s

210    or 7 s, and delay length and sample contrast (*high*; *low*; *null*) were fully crossed. For both

211    Experiments 1A and 1B, only trials with a 7-s delay were included in the behavioral analyses to

212    better match the duration of the fMRI task.

213

214    **Behavioral analyses**

215    Behavioral performance was assessed in two ways. Within-trial recall error was

216    calculated for *high* and *low* trials as the angular difference between the sample and response

217    orientations, for each condition separately. Differences between conditions were evaluated by

218    paired t-tests. Serial bias on response from the previous trial was calculated for all three

219    conditions. This was done by calculating the difference between the current and previous

220    response, and grouping the difference values into nine 20º-wide bins. To test whether the number

221    of trials differed between bins, we performed a $\chi^2$ goodness of fit test on each condition.

222

223    **fMRI methods**

224    *Data acquisition*

225    Whole-brain images were acquired with a 3 Tesla GE MR scanner (Discovery MR750; GE

226    Healthcare, Chicago, IL) with a 32-channel head coil at the Lane Neuroimaging Laboratory at

227    the University of Wisconsin–Madison HealthEmotions Research Institute (Department of

228    Psychiatry). Functional images were acquired with a gradient-echo echo-planar sequence (2 s

229    repetition time (TR), 22 ms echo time (TE), 60° flip angle) within a 64 × 64 matrix (42 axial

230    slices, 3 mm isotropic). A high-resolution T1 image was also acquired for each session with a

231    fast spoiled gradient-recalled-echo sequence (8.2 ms TR, 3.2 ms TE, 12° flip angle, 176 axial

232    slices, 256 × 256 in-plane, 1.0 mm isotropic).

233

234    *Preprocessing*

235    Functional MRI data were preprocessed using AFNI (http://afni.nimh.nih.gov) (Cox, 1996). The

236    first four volumes of each functional run were removed. The data were then registered to the first

237    volume of the first run within each scan session, and then to the T1 volume of the same session.

238    Data from the second session were further registered to the T1 volume of the first scanning

10

239    session. The data were then motion corrected, detrended (linear, quadratic, cubic), converted to

240    percent signal change. Data for subsequent general linear model (GLM) analyses were further

241    spatially smoothed with a 4-mm FWHM Gaussian kernel. Data for multivariate and univariate

242    time course analyses were zscored within each run.

243

244    *Univariate analyses*

245    Task-related changes in activity were identified with a mass-univariate GLM implemented in

246    AFNI, with sample, delay and probe epochs of the task modeled with boxcars (0.5 s, 8.5 s or 9.5

247    s depending on the participant, and 4 s, respectively), each convolved with a canonical

248    hemodynamic response function. Six nuisance regressors were also included to account for head

249    motion artifacts in the six dimensions of rigid body motion.

250    Percent signal change in BOLD activity relative to baseline was calculated for each time point

251    during the working memory task, baseline was chosen as the average BOLD activity of the first

252    TR of each trial. BOLD signal change was averaged across trials within each condition, and

253    across all voxels within each region of interest (ROI; see below).

254    Statistical significance of BOLD activity against baseline was assessed using two-tailed, one-

255    sample t-tests against 0, and the resultant $p$ values were corrected across time points using FDR

256    (False Discovery Rate) (Benjamini & Hochberg, 1995). Statistical difference of BOLD activity

257    between conditions at each time point was assessed using two-tailed paired t-tests, with FDR

258    correction applied across time points and comparisons.

259

260    *Region of interest (ROI) definition*

261    We created subject-specific anatomical ROIs by warping masks from the probabilistic atlas of

262    Wang and colleagues (2015) to each subject's structural scan in their native space. Early visual

263    anatomical ROIs were created by merging the masks for unilateral V1, V2, and V3 within and

264    between hemispheres. IPS anatomical ROIs were created by merging the masks for unilateral

265    regions IPS0-5 within and between hemispheres. For the *Early Visual Cortex* functionally

266    defined ROI, we identified the 500 voxels displaying the strongest loading on the contrast

267    [sample - baseline], collapsing over all three conditions. For the *IPS* functionally defined ROI,

268    we identified the 500 voxels displaying the strongest loading on the contrast [delay - baseline],

269    collapsing over all three conditions. For completeness, an alternate "*Sample IPS* ROI" was also

270    defined as the 500 voxels in this anatomical region displaying the strongest loading on the

271    contrast [sample - baseline].

272

273    *Multivariate inverted encoding modeling*

274    All inverted encoding modeling (IEM) analyses were performed using custom functions in

275    MATLAB. The IEM assumes that the responses of each voxel can be characterized by a small

276    number of hypothesized tuning channels. Following previous work, the number of orientation

277    tuning channels was set to nine (20º apart, equally spaced), and the idealized feature tuning curve

278    of each channel to a specific orientation $\theta$ was defined as a half-wave-rectified sinusoid raised to

279    the eighth power (FWHM = 0.82 rad):

280    $f(\theta) = \cos(\theta - c)^8$

281    Where $c$ was the center of the channel.

282          We then computed the weight matrix ($W$, $v \times k$, $v$: the number of voxels; $k$: the number of

283    channels) that projects the hypothesized channel responses ($C_1$, $k \times n$, $n$: the number of trials) to

284    actual measured fMRI signals in the training dataset ($B_1$, $v \times n$), and extracted the estimated

285    channel responses ($\hat{C}_2$, $k \times n$) for the test dataset ($B_2$, $v \times n$) using this weight matrix.

12

286    The relationship between the training dataset ($B_1$) and the channel responses ($C_1$) was

287    characterized by:

288    $B_1 = W C_1$

289    Therefore, the least-squared estimate of the weight matrix ($\widehat{W}$) was calculated using

290    linear regression:

291    $\widehat{W} = B_1 C_1^T (C_1 C_1^T)^{-1}$

292    The channel responses ($\hat{C}_2$) for the test dataset ($B_2$) was then estimated using the weight

293    matrix ($\widehat{W}$):

294    $\hat{C}_2 = (\widehat{W}^T \widehat{W})^{-1} \widehat{W}^T B_2$

295    Because orientations in the current study were randomly selected from the 1º - 180º

296    orientation space (in steps of 1º), we did not pick a fixed set of channel centers, as is often done

297    (Yu & Shim, 2017; Yu, Teng, et al., 2020). Instead, following Rademaker et al. (2019) we first

298    picked a set of equally spaced channel centers (e.g., 0º, 20º, 40º, 60º, 80º, 100º, 120º, 140º, 160º),

299    conducted the analysis as described above, and then shifted the channel centers by 1º and

300    repeated the analysis. The procedure was repeated 20 times, such that all 180 orientations from

301    1º to 180º in 1º step served as channel centers. We then combined estimated channel responses

302    from all iterations of these analyses to create responses of 180 orientation channels. The result,

303    for any given orientation, can be considered a reconstruction of the model's estimate of the

304    neural representation of that orientation. This procedure ensured that our channel estimates were

305    not biased by any specific channel centers. All channel responses were then centered on a

306    common center (0º on the x-axis) and averaged for visualization and for statistical comparisons.

307

308    *Hypothesis testing*

309    *Analysis plan.* If a participant is not aware of the fact that a considerable proportion of

310    trials will feature *null* samples that contain no orientation information, we assume that on *null*

311    trials they will generate an orientation for response at some point prior to the onset of the

312    response dial. Furthermore, because the initial orientation of the dial cannot be predicted prior to

313    its onset, we assume that this response plan will not be kinematic (e.g., how many degrees they

314    plan to turn the dial), but rather will be the representation of the orientation that the participant

315    plans to produce at the end of trial. To validate this assumption, our first analysis would be to

316    train an IEM using the orientation of the response on that trial (response-based IEM). Successful

317    reconstruction of orientation with this IEM at time points preceding the response (i.e., during the

318    delay period) would mean that participants were indeed representing the orientation of their

319    planned response during those earlier time points (response-based neural code)."

320    Assuming success of this first analysis, the next step would be to determine whether a

321    common response-based neural code was employed across conditions. This would be done by

322    applying the response-based IEM from one trial type (e.g., *high*) to data from the other two trial

323    types (e.g., *low* and *null*). We anticipated three possible outcomes: If reconstruction in a tested

324    condition was significantly positive, and did not differ from that in the training condition, this

325    would reflect "full generalization"; if reconstruction in a tested condition was significantly

326    positive, but was also significantly lower than that in the training condition, this would reflect

327    "partial generalization"; and if reconstruction in a tested condition was not significant, this would

328    reflect "failed generalization". These results would be interpreted as evidence for a fully shared

329    neural code, for a partly shared neural code, or as a failure to find evidence for a shared neural

330    code, respectively.

331    Finally, because IEM relies on specific hypotheses of orientation channels, we would also

332    perform a model-free analysis, multidimensional scaling (MDS), to see if this alternative

333    approach would support conclusions comparable to those suggested by the IEM analyses.

334    *Operationalizing hypothesis tests.* To investigate the codes supporting the representation

335    of orientation in the different conditions (*high*; *low*; *null*), we trained two IEMs: a response-

336    based IEM labeled according to the orientation of the response on each trial, and a sample-based

337    IEM labeled according to the sample orientation on each trial. Note that the response-based IEM

338    would be the focus of our analyses, and results from the sample-based IEM in the *null* condition

339    would not be interpretable on its own, but would serve as controls for comparing with the results

340    from *high* and *low* conditions. IEMs were trained and tested using a leave-one-run-out cross-

341    validation procedure, for each condition, time point (or average of time points, e.g., average of

342    time points 8 - 10 s for delay period), and ROI separately. To compare response-based neural

343    codes across conditions, we also used a leave-one-run-out procedure, training the response-based

344    IEM on data from one condition, and testing the IEM on the data from all three conditions,

345    including the training condition (which would yield the same result as the first analysis) and the

346    two other conditions. This procedure was performed for each condition, time point (or average of

347    time points), and ROI separately.

348    We also trained several complementary IEMs for testing alternative explanations for the

349    results. First, we trained a mixed IEM using a balanced set of trials from each condition (*high*;

350    *low*; *null*), and tested this IEM on the same balanced set of trials from each condition separately.

351    The purpose of this IEM would be to avoid potential concerns with differences in SNR across

352    IEMs (Liu, Cable, & Gardner, 2018; Sprague et al., 2018). Second, to examine the influence of

353    previous-trial information on the reconstruction of current-trial information, we trained response-

354  based IEMs using response labels from the previous trial, or trained sample-based and response-

355  based IEMs while excluding trials with similar response to that of the previous trial.

356      To characterize the strength of each IEM reconstruction, we collapsed over the channel

357  responses on both sides of the common center, averaged them, then calculated the slope of each

358  collapsed reconstruction using linear regression (Foster, Bsales, Jaffe, & Awh, 2017; Samaha,

359  Sprague, & Postle, 2016). A larger positive slope indicates stronger positive representation. We

360  used a bootstrapping procedure (Ester et al., 2015; Yu, Teng, et al., 2020) to characterize the

361  significance of the slopes. For each combination of factors (IEM, condition, time point, or ROI),

362  seventeen orientation reconstructions were randomly sampled with replacement from the pool of

363  seventeen participants and averaged. This procedure was repeated 10000 times, resulting in

364  10000 average orientation reconstructions, and correspondingly 10000 slopes. The probability of

365  obtaining a negative slope among the 10000 slopes was counted as the one-tailed $p$-value of the

366  slope. To characterize the difference between the slopes of two IEM reconstructions, we first

367  calculated the difference between two bootstrapped slopes 10000 times, which generated 10000

368  slope differences. The significance of the slope difference was then calculated using the same

369  one-tailed method as above. All $p$-values were corrected for multiple comparisons using the FDR

370  method, across IEMs (sample-based, response-based), conditions (*high*, *low*, or *null*), and time

371  points.

372      We also assessed evidence for differences between the slopes of delay-period response-

373  based reconstructions with Bayes Factors (BF), which support evaluation of the amount of

374  evidence for one hypothesis ($H_1$) against the null hypothesis ($H_0$). $H_1$ referred to a positive

375  reconstruction and $H_0$ referred to a failed reconstruction (i.e., a slope no larger than 0). For

376  comparison between the slopes of two reconstructions, $H_1$ referred to the slopes being different

377  and $H_0$ referred to the absence of evidence for a difference. As an example, a $BF_{10}$ of 3 would

378    indicate that $H_1$ is three times more likely than $H_0$, whereas a $BF_{10}$ of 0.33 would indicate that $H_0$

379    is three times more likely than $H_1$. All BF analyses were conducted using the JASP software

380    (Love et al., 2019).

381

382    *Multidimensional scaling*. For each ROI and each trial epoch we categorized all response

383    orientations into four bins (0-45º, 45-90º, 90-135º, 135-180º). Trial number for each condition

384    was matched by subsampling data from the *high* and *null* conditions to match the number of

385    trials in the *low* condition. The Euclidean distances between orientation bins and conditions were

386    then computed using the covariance matrix calculated from the subsampled data. This

387    subsampling procedure was repeated for 1000 times and averaged. Distances were averaged

388    across participants, and multidimensional scaling was performed on the distance matrix using the

389    "cmdscale" function in MATLAB.

390

391    **Results**

392    **Behavior**

393    *Experiment 1A*

394    Participants' 75% contrast discrimination threshold for recall of orientation against a noise

395    background ranged between 4-6%, with a mean of 5.0% and a standard deviation of 0.6%. For

396    delayed recall of the orientation of a sample grating, the average recall error for *high* contrast

397    (60%) samples ($9.0º \pm 1.5º$) was significantly lower than for the *threshold* contrast samples

398    ($17.4º \pm 4.4º$), $t(12) = 5.95$, $p < 0.001$.

399

400    *Experiment 1B*

401   Delayed recall of orientation did not differ between *high* contrast (60%; 10.2º ± 1.9º) and *low*

402   contrast (10%; 10.6º ± 2.7º) conditions, $t(15) = 0.64$, $p = 0.530$ (Figure 1). The fact that average

403   recall error did not differ between the *high* and *low* trials established the fact, critical for the logic

404   of Experiment 2, that *low* and *high* samples were comparably visible to participants. This, plus

405   the marked difference between performance at these two levels of contrast versus performance at

406   75% threshold (Experiment 1A) indicated that neither *high* nor *low* contrast trials were likely to

407   produce trials in which the sample grating was not visible to the participant (in contrast to some

408   of the *threshold* trials in Experiment 1A).

409

410   *Experiment 2*

411   Consistent with Experiment 1B, Recall error during scanning did not differ between the *high*

412   (11.4º ± 4.0º) and *low* (11.7º ± 4.8º) trials, $t(16) = 0.58$, $p = 0.567$.

413        Although recall error could not be calculated for *null* trials, the results from several

414   analyses suggested that participants did not treat *null* trials different from trials on which a

415   sample grating was visible. First, angular difference between the starting position of the response

416   needle and the recalled orientation did not differ between *high*, *low*, and *null* trials (42.3º ± 2.4º,

417   41.7º ± 3.9º, 41.8º ± 5.7º, respectively; all $t$s < 0.85, $p$s > 0.408), suggesting that the three

418   conditions were comparable in terms of effort during recall. Second, although sample orientation

419   on each trial was randomly chosen and the distribution of sample orientations was uniform (i.e.,

420   there was an equal proportion of cardinal and oblique orientations), plotting the distribution of

421   participants' raw responses showed biased responses towards oblique orientations (relative to

422   cardinal orientations) for all three trial types (Figure 2). This indicates that trials of all types were

423   influenced to a similar extent by a systematic bias, perhaps from one or more stimulus-

424   nonspecific factors such as prior knowledge (Panichello et al., 2019; Yu, Panichello, Cai, Postle,

425 & Buschman, 2020). In sum, *null* and *high*/*low* trials were well matched in terms of procedural

426 details, and the only difference between conditions was the availability of external orientation

427 information. Therefore, any orientation information observed in the *null* trials could only be

428 internally generated.

429

430 <----- insert Figure 2 about here ----->

431

432 **fMRI**

433 *Time course of BOLD activity*

434 All analyses were carried out at the level of the *Early Visual Cortex* ROI and the *IPS* ROI. In

435 both regions, a conventional time course of BOLD activity change was observed for all three

436 conditions (Figure 3): sample-evoked activity reached its peak at around 4-6 s after trial onset;

437 delay-period activity reached its trough at around 8-10 s; and response-evoked activity reached

438 its peak at around 14-16 s. Time points 8-10 s were subsequently used to operationalize "late

439 delay-period" activity. In early visual cortex, activity in *null* trials was slightly lower than that in

440 *high* and *low* trials during sample and early delay epochs (2-8 s; all $p$s < 0.023), but not at 10 s

441 (both $p$s > 0.167) nor during the response epoch (12 s and after; all $p$s > 0.342). In IPS, in

442 contrast, *null* activity was lower during the sample (2-4 s; all $p$s < 0.005) and response epochs

443 (12-18 s, all $p$s < 0.040 except for 12 s between *high* and *null*: $p$ = 0.073), but not during the

444 delay (6-10 s, all $p$s > 0.132).

445

446 <----- insert Figure 3 about here ----->

447

448 *Inverted encoding modeling*

449     *Early visual cortex.* To assess the time course of neural representations of orientation, for

450     each trial type, we applied a sample-based IEM (i.e., trained on the sample label) and a response-

451     based IEM (trained on the response label) to every time point of the trial. For *high* and *low* trials,

452     reconstruction with the sample-based IEM was significant beginning at 4 s after sample onset

453     and sustained for the remainder of the trial (all $p$s < 0.001). Similarly, reconstruction with the

454     response-based IEM were significant for the duration of trial, beginning at 4 s for *high* trials and

455     at 2 s for *low* trials (all $p$s < 0.040). Sample-based and response-based reconstructions did not

456     differ at any time point, for either of these two conditions (all $p$s > 0.157). These results validated

457     the approach of using participants' responses as an estimate of the orientation that they

458     represented earlier in the trial, prior to the response.

459     Turning next to *null* trials, reconstruction with sample-based IEMs did not achieve

460     statistical significance except for two isolated time points: 2 s ($p$ = 0.017) and 16 s ($p$ = 0.036),

461     probably due to statistical noise. Note that these null results amounted to confirmation of a sanity

462     check, because the labels used to train the sample-based IEM did not correspond to what subjects

463     were presented on these trials. Reconstructions with response-based IEMs, were significant

464     beginning with 6 s and for the duration of trial (all $p$s < 0.020; Figure 4A). Critically, these

465     response-based reconstructions were significantly different from the sample-based

466     reconstructions for 6-8 s and from 12 s onwards (green asterisks; all $p$s < 0.012), suggesting that

467     robust orientation representations specific to the response started from 6 s after trial onset. This

468     indicates that, beginning relatively early in the trial, participants generated and maintained a

469     representation with exclusively internally derived information.

470     *Intraparietal sulcus.* In IPS, results were generally comparable to those from early visual

471     cortex, albeit weaker in magnitude. When focusing on the late-delay period (Figure 4D), sample

472     and response reconstructions were significant in all conditions (all $p$s < 0.037), except for the

473   sample reconstruction in the null condition ($p = 0.259$). Time point-by-time point reconstructions

474   were also qualitatively similar to early visual cortex (Figure 4C): on *high* trials sample and

475   response reconstructions emerged during the sample period and were sustained throughout the

476   trial, as were sample reconstructions on *low* trials (all $p$s < 0.041). Response reconstructions

477   were smaller in slope on *low* trials, and, with the exception of a single time point (6 s after trial

478   onset, $p = 0.007$), did not survive correction for multiple comparisons during the delay. Note that

479   the lower number of trials for the *low* condition might have been responsible for the lack of

480   significance here. Indeed, robust reconstruction of orientation was observed for *low* trials when

481   averaging across time points in the delay period (Figure 4D).

482        Turning to *null* trials, reconstructions with sample-based IEMs only achieved statistical

483   significance at 2 s, a result probably be due to statistical noise. Reconstruction with response-

484   trained IEMs, however, was significant for all time points beginning with 4 s (all $p$s < 0.028),

485   with the exception of 10 s of the delay period ($p = 0.076$).

486        We also carried out these analyses in the *Sample IPS* ROI (IPS ROI defined using

487   sample-period activity), and the results (not shown) were qualitatively similar to those in the

488   *Delay IPS* ROI.

489

490                          <----- insert Figure 4 about here ----->

491

492        One possible concern about the finding of principal theoretical interest from these

493   analyses -- the reconstruction of response-related stimulus information from the delay period of

494   *null* trials (Figure 4) -- is that this might reflect "spillover" of information processed during the

495   previous trial, rather than evidence for genuinely internally generated stimulus representations.

496   Additional analyses carried out to assess this alternative possibility ruled it out as a major

21

497     concern, and these are presented at the end of the *Results* section (see *Secondary analyses to*

498     *assess the influence of the previous trial on response-based IEMs*).

499

500     *Comparison of neural codes across high, low, and null trials*

501     Having established robust measurements for internally generated neural representations of

502     orientation, we next sought to examine the nature of these representations. Specifically, because

503     representations on *null* trials were purely internally generated, whereas representations on *high*

504     and *low* trials reflected influences from both external and internal sources, we tested whether the

505     representations maintained during these different trial types recruited a common neural code, in

506     keeping with previous demonstrations of a shared neural code between working memory and

507     perception (Harrison & Tong, 2009), between working memory and attention (Yu & Shim,

508     2019), and between working memory and imagery (Albers et al., 2013). To this end we trained

509     IEMs on one condition, and tested it on all three conditions (see *Methods*). Note that only

510     response-based IEMs were recruited for this purpose. For these analyses, we emphasized the

511     results from the late delay period (8-10 s after trial onset; also see Figure 5 for results for the full

512     time courses). Here we also employed Bayes Factors (BF) to assess the amount of evidence in

513     generalization. A BF of larger than 3 or smaller than 1/3 can be considered substantial evidence

514     supporting or rejecting the hypothesis.

515

516                    <----- insert Figure 5 about here ----->

517

518         In early visual cortex, we successfully reconstructed orientation from the late-delay

519     period of *low* trials with IEMs trained on *high* trials ($p < 0.001$, $BF_{10} = 280.3$), and of *high* trials

520     with IEMs trained on *low* trials ($p < 0.001$, $BF_{10} = 79.5$). Furthermore, these results demonstrated

521  full generalization: reconstructions for *high* and *low* trials with the *high*-trained IEM did not

522  differ from each other ($p = 0.552$, $BF_{10} = 0.5$); nor did reconstructions for *high* and *low* trials

523  with the *low*-trained IEM ($p = 0.477$, $BF_{10} = 0.8$; Figure 6A and 6B). When comparing each of

524  these visible trial types with *null* trials, in contrast, cross-condition generalization was

525  asymmetric: For *high* trials, although the IEM trained on *high* trials failed to generalize to *null*

526  trials ($p = 0.135$, $BF_{10} = 0.7$; Figure 6A), the IEM trained on *null* trials did successfully

527  reconstruct orientation on *high* trials ($p = 0.0014$, $BF_{10} = 6.0$), and reconstructions for *high* and

528  *null* trials with the *null*-trained IEM did not differ from each other ($p = 0.552$, $BF_{10} = 0.2$; Figure

529  6C). For *low* trials, on one hand, the IEM trained on *null* trials successfully reconstructed

530  orientation on *low* trials ($p = 0.0013$, $BF_{10} = 8.8$), and reconstructions for *low* and *null* trials with

531  the *null*-trained IEM did not differ from each other ($p = 0.552$, $BF_{10} = 0.2$; Figure 6C); On the

532  other hand, the IEM trained on *low* trials did generalize to *null* trials ($p = 0.0052$, $BF_{10} = 5.3$),

533  although the slope of this reconstruction was lower than that on *low* trials with the *low*-trained

534  IEM ($p = 0.030$, $BF_{10} = 28.4$; Figure 6B), suggesting only partial generalization from *low* to *null*

535  trials.

536      In IPS, although response-based neural codes were also fully generalizable between *high*

537  and *low* trials (train *high*-test *low*, $p = 0.007$, $BF_{10} = 6.2$; train *low*-test *high*, $p = 0.006$, $BF_{10} =$

538  6.9; train *high*-test *high* vs. train *high*-test *low*, $p = 0.732$, $BF_{10} = 0.2$; train *low*-test *low* vs. train

539  *low*-test *high*, $p = 0.497$, $BF_{10} = 0.3$; Figure 6D and 6E), there was no evidence for cross-

540  generalization from *null* trials to *high* or *low* trials (train *null*-test *high*, $p = 0.215$, $BF_{10} = 0.5$;

541  train *null*-test *low*, $p = 0.061$, $BF_{10} = 1.6$; Figure 6F), nor from *high* or *low* trials to *null* trials

542  (train *high*-test *null*, $p = 0.252$, $BF_{10} = 0.4$; and train *low*-test *null*, $p = 0.187$, $BF_{10} = 0.6$; Figure

543  6D and 6E).

544

545                                                  <----- insert Figure 6 about here ----->

546

547        Although cross-generalization is a common approach for assessing commonality of

548    neural codes (Albers et al., 2013; Rademaker et al., 2019; Yu & Shim, 2019), interpreting

549    failures to generalize can be complicated by technical considerations arising from training the

550    IEM on the same versus on different datasets (Liu, Cable, & Gardner, 2018; Sprague et al.,

551    2018). Therefore, we repeated these analyses but with a single IEM trained on a balanced set of

552    trials drawn in equal number from *high*, *low*, and *null* trials. Results with this mixed IEM were

553    complementary to the cross-generalization analyses: in both early visual cortex and IPS, the

554    mixed IEM generated successful reconstructions of orientation from *high* and *low* trials (4-10 s:

555    all $p$s < 0.006), but failed on *null* trials (4-10 s: all $p$s > 0.140; Figure 7).

556

557                                                     <----- insert Figure 7 about here ----->

558

559    *Model-free analyses*

560    Lastly, to determine whether a difference between *null* and *high*/*low* trials would be observed

561    when no model-based approach was applied to the data, we compared the representational

562    distances between conditions using MDS. MDS analyses were performed for the sample (4-6 s

563    after trial onset), delay (8-10 s after trial onset), and response (14-16 s after trial onset) epochs of

564    the working memory task, separately for early visual cortex and for IPS. For visualization

565    purposes, response orientations were grouped into four 45º-wide bins.

566        In early visual cortex, during the sample epoch, the three conditions were discriminable

567    along Dimension 1, confirming that differences in stimulus contrast influenced sensory

568    processing (Figure 8A). During the delay period the distance between the *high* and *low*

569   conditions decreased, such that the two now overlapped along Dimension 1, while the *null*

570   condition remained separated from the other two. This suggested that, as stimulus-driven

571   influences diminished, trials that relied exclusively on internally derived information remained

572   distinct. This discriminative element carried on into the response period, along Dimensions 2 and

573   3, despite the fact that participants performed the same type of motor response on every trial. In

574   IPS, a similar discriminative pattern was also observed between conditions (Figure 8B). Thus, in

575   both brain areas, patterns of activity on *null* trials were distinct from those on *high*/*low* trials in

576   multidimensional representational space. The fact that this was true for all epochs of the trial

577   suggests that this separability was not simply a result of perceptual differences between memory

578   samples.

579

580                       <----- insert Figure 8 about here ----->

581

582   *Secondary analyses to assess the influence of the previous trial on response-based IEMs*

583   Recent perceptual history can bias behavior on the current trials (Fischer & Whitney, 2014),

584   including during working memory tasks (Barbosa et al., 2020; Samaha, Switzky, & Postle,

585   2019), and it has been shown that the no-longer-relevant content of the previous trial can be

586   decoded from electroencephalography (EEG) signals recorded during a visual working memory

587   task (Bae & Luck, 2019). Consequently, we carried out a series of analyses to assess whether the

588   response-related reconstructions from *null* trials (Figure 4), rather than reflecting internally

589   generated stimulus representations, might instead be due to "spillover" of information processed

590   during the previous trial. We tested this possibility with two approaches. First, we examined if

591   the response of the previous trial could be reconstructed from patterns of activity of the current

592   trial in the current data. In early visual cortex we found that the response of the previous trial

593   could indeed be reconstructed in all conditions, especially during the earlier portion of the trial,

594   (all $p$s < 0.031; Figure 9). However, because above-baseline-level reconstruction of the previous-

595   trial response was present at the very beginning of the trial (i.e., 0 s), and reconstruction of the

596   current-trial response did not emerge until 6 s after trial onset, we believe it was unlikely that

597   these two sets of results reflected the same piece of information. Furthermore, in IPS,

598   reconstruction of the previous trial's response was almost absent, with the exception of three

599   isolated time points across all three conditions (all $p$s < 0.03). This effect alone again cannot

600   explain the sustained reconstructions of the response orientation on *null* trials.

601

602                           <----- insert Figure 9 about here ----->

603

604        A second approach to assess the possible influence of information from previous trials on

605   response-related reconstructions from *null* trials was to redo the analyses after removing the

606   trials for which the response was most similar to the response on the previous trial. We did this

607   by first calculating the difference between each trial's response and the response on the previous

608   trial, for all three conditions, and grouping the trials by difference values into nine 20º-wide bins.

609   For *high* and *null* trials, the distribution of the differences was not uniform ($\chi^2(8)$ = 19.5 and

610   81.9, $p$ = 0.012 and $p$ < 0.001, respectively; Figure 10), suggesting a potential influence from

611   previous trials on the performance of the current trial. Next, for *null* trials, we removed the

612   influence from the responses that were closest to the previous response (difference < 10º; bins

613   highlighted in red in Figure 10) by excluding trials that belonged to this bin and repeating the

614   IEM analyses on the remaining trials. Significant response reconstructions were still observed in

615   this subset of *null* condition trials (Figure 11), increasing our confidence that the representation

26

616     of response-related orientation information on *null* trials cannot be simply explained as

617     reactivation of perceptual history from the previous trial.

618

619                     <----- insert Figures 10 and 11 about here ----->

620

621         Finally, we examined whether the potency of the spillover effect varied with sample type,

622     by sorting every *high* trial as a function of whether it was preceded by a *high*, *low*, or *null* trial.

623     Results (not shown) indicated that the spillover effect was comparable for each trial type, and

624     that the time course of each mimicked the pattern seen in Figure 9.

625

626     **Discussion**

627         The human brain processes massive amounts of information every day, from both

628     external and internal sources. To explore how internally-generated information is represented in

629     the brain during working memory, we incorporated a null-sample condition into a delayed-recall

630     task. First, we demonstrated that, after the presentation of a null sample, participants generated a

631     neural representation that corresponded to the response that they would make at the end of the

632     trial, confirming that our procedure was successful at producing internally generated working-

633     memory representations. Next, we assessed cross-generalization of the neural representation of

634     orientations between conditions, and observed an asymmetric pattern in early visual cortex:

635     IEMs trained on data from *null* trials generalized fully to data from visible-sample trials, but the

636     converse was not true. This suggested some difference in the processing of internally generated

637     representations versus conventional working memory representations that receive influences

638     from both external and internal sources. This difference in neural codes was also evident when

639     the data were projected into multidimensional representational space: the patterns of activity for

640   *high* and *low* trials were clearly segregated from *null* trials in both early visual cortex and IPS.

641   Therefore, stimulus information that is derived from an external source is represented differently

642   than stimulus information that is generated internally.

643   Our findings might seem inconsistent with previous work that has demonstrated a shared

644   neural code between visual working memory and visual imagery in early visual cortex (Albers et

645   al., 2013). However, because visual imagery tasks often involve elements such as mental

646   rotations (Albers et al., 2013) or retrocueing manipulations (Dijkstra et al., 2018), they typically

647   refer overtly to previously presented (i.e., externally originated) information, and this may

648   explain why similar neural codes are recruited by these two classes of task. It had thus remained

649   unclear whether "purely" internally-derived representations also share the same neural code as

650   "conventional" working memory representations. The present results -- indicating that the

651   representation of orientation in early visual cortex fully generalizes from the *null* to the *high* and

652   *low* conditions, but not in the other direction – suggest that all three conditions share the same

653   purely internally generated neural codes, and that conventional working memory representations

654   contain one or more additional dimensions that are lacking from "purely" internally originated

655   visual representations. The additional dimension(s) are likely related to processes that are

656   involved in the initial processing of externally presented information.

657   The differences between working memory and internally originated imagery were also

658   preserved in IPS, where we found that *null* and *high*/*low* trials did not generalize in either

659   direction, although the effects were generally weaker compared with those in early visual cortex.

660   These results are in line with previous work demonstrating failures to find evidence -- in higher-

661   order parietal and/or frontal cortex -- for generalization of neural codes between working

662   memory and visual perception (Rademaker et al., 2019), attention (Yu & Shim, 2019), and

663   imagery (Albers et al., 2013).

664    What is the nature of the internally generated representations observed in the delay period

665    of the *null* condition in the current study? One possibility is a preparatory motor code, similar to

666    what has been demonstrated for visual working memory for orientation on a task that allowed for

667    concurrent selection of visual and motor responses (van Ede, Chekroud, Stokes, & Nobre, 2019).

668    If so this would need to be a highly abstract code, akin to an intention, because the starting

669    position of the probe in our experiment was randomized from trial-to-trial, and so participants

670    would not have been able to plan their specific motor response prior to the onset of the response

671    wheel. Another possibility is that they reflected internally generated representations of

672    participants' best guess of the orientation of the sample. This would be consistent with the fact

673    that the time course of the representation of orientation developed later in time in the *null*

674    condition relative to the *high* and *low* conditions, especially in early visual cortex. Similarly, it

675    has been observed that representations of visual imagery develop later in time than do

676    representations associated with visual perception (Dijkstra et al., 2018) and with visual working

677    memory (Albers et al., 2013). It is likely that internally generated representations are influenced

678    by many stimulus-nonspecific factors, such as prior knowledge (Bae et al., 2015; Bae et al.,

679    2014; Panichello et al., 2019; Yu, Panichello, et al., 2020) and recent history (Bae & Luck, 2019;

680    Fischer & Whitney, 2014), and that these stimulus-nonspecific factors may serve, in part, as

681    differentiating factors in the coding of externally driven versus internally generated information.

682    Indeed, we did observe influences from the previous trial in the current experiment, although

683    spillover from the previous trial alone cannot explain the sustained, robust representations of the

684    response in the *null* condition. It should be possible to use the "null-sample" paradigm in

685    combination with other visual tasks to better understand the nature of internally generated visual

686    representations. For example, it would be interesting to compare internally generated

687    representations directly with the codes that support visual perception. It would also be interesting

688    to include confidence ratings in future tasks to better understand the subjective experience of the

689    null task. Finally, by combining the paradigm with ultra-high field fMRI, one would be able to

690    investigate whether there exist layer-specific representations for purely internally generated

691    representations.

692         Our results, together with previous work (Albers et al., 2013; Harrison & Tong, 2009;

693    Rademaker et al., 2019; Yu & Shim, 2019), suggest a potential mechanism for how the brain

694    processes information originating from different sources. Early visual cortex represents stimulus

695    properties with a common neural code that is insensitive to behavioral/cognitive context, such

696    that the same neural code is shared between visual perception, attention, and working memory,

697    consistent with the sensorimotor recruitment hypothesis. However, early visual cortex also

698    registers the source of origination of this information, such that externally originated and

699    internally originated representations can be differentiated. This distinction between externally

700    and internally originated representations was also observed in a higher-order cortical area, IPS,

701    although perhaps with a slightly different pattern. These signals may underlie the neural basis for

702    how the brain differentiates and maintains signals from different sources.

703

704    **References**
705
706    Albers, A. M., Kok, P., Toni, I., Dijkerman, H. C., & de Lange, F. P. (2013). Shared
707         representations for working memory and mental imagery in early visual cortex. *Curr
708         Biol, 23*(15), 1427-1431. doi:10.1016/j.cub.2013.05.065
709    Awh, E., & Jonides, J. (2001). Overlapping mechanisms of attention and spatial working
710         memory. *Trends Cogn Sci, 5*(3), 119-126.
711    Baddeley, A. (2003). Working memory: looking back and looking forward. *Nat Rev Neurosci,
712         4*(10), 829-839. doi:10.1038/nrn1201
713    Bae, G. Y., & Luck, S. J. (2019). Reactivation of Previous Experiences in a Working Memory
714         Task. *Psychol Sci, 30*(4), 587-595. doi:10.1177/0956797619830398
715    Bae, G. Y., Olkkonen, M., Allred, S. R., & Flombaum, J. I. (2015). Why some colors appear
716         more memorable than others: A model combining categories and particulars in color
717         working memory. *J Exp Psychol Gen, 144*(4), 744-763. doi:10.1037/xge0000076

718  Bae, G. Y., Olkkonen, M., Allred, S. R., Wilson, C., & Flombaum, J. I. (2014). Stimulus-specific
719      variability in color working memory with delayed estimation. *Journal of Vision, 14*(4), 7-
720      7. doi:10.1167/14.4.7
721  Barbosa, J., Stein, H., Martinez, R. L., Galan-Gadea, A., Li, S., Dalmau, J., . . . Compte, A.
722      (2020). Interplay between persistent activity and activity-silent dynamics in the prefrontal
723      cortex underlies serial biases in working memory. *Nat Neurosci*. doi:10.1038/s41593-
724      020-0644-4
725  Benjamini, Y., & Hochberg, Y. (1995). Controlling the False Discovery Rate - a Practical and
726      Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B-
727      Statistical Methodology, 57*(1), 289-300.
728  Bettencourt, K. C., & Xu, Y. (2016). Decoding the content of visual short-term memory under
729      distraction in occipital and parietal areas. *Nat Neurosci, 19*(1), 150-157.
730      doi:10.1038/nn.4174
731  Brainard, D. H. (1997). The Psychophysics Toolbox. *Spat Vis, 10*(4), 433-436.
732  Cox, R. W. (1996). AFNI: software for analysis and visualization of functional magnetic
733      resonance neuroimages. *Comput Biomed Res, 29*(3), 162-173.
734  D'Esposito, M., & Postle, B. R. (2015). The cognitive neuroscience of working memory. *Annu
735      Rev Psychol, 66*, 115-142. doi:10.1146/annurev-psych-010814-015031
736  Dijkstra, N., Mostert, P., Lange, F. P., Bosch, S., & van Gerven, M. A. (2018). Differential
737      temporal dynamics during visual imagery and perception. *Elife, 7*.
738      doi:10.7554/eLife.33904
739  Engle, R. W., Tuholski, S. W., Laughlin, J. E., & Conway, A. R. (1999). Working memory,
740      short-term memory, and general fluid intelligence: a latent-variable approach. *J Exp
741      Psychol Gen, 128*(3), 309-331.
742  Ester, E. F., Sprague, T. C., & Serences, J. T. (2015). Parietal and Frontal Cortex Encode
743      Stimulus-Specific Mnemonic Representations during Visual Working Memory. *Neuron,
744      87*(4), 893-905. doi:10.1016/j.neuron.2015.07.013
745  Fischer, J., & Whitney, D. (2014). Serial dependence in visual perception. *Nat Neurosci, 17*(5),
746      738-743. doi:10.1038/nn.3689
747  Foster, J. J., Bsales, E. M., Jaffe, R. J., & Awh, E. (2017). Alpha-Band Activity Reveals
748      Spontaneous Representations of Spatial Position in Visual Working Memory. *Curr Biol,
749      27*(20), 3216-3223 e3216. doi:10.1016/j.cub.2017.09.031
750  Gilbert, C. D., & Li, W. (2013). Top-down influences on visual processing. *Nat Rev Neurosci,
751      14*(5), 350-363. doi:10.1038/nrn3476
752  Harrison, S. A., & Tong, F. (2009). Decoding reveals the contents of visual working memory in
753      early visual areas. *Nature, 458*(7238), 632-635. doi:10.1038/nature07832
754  Lawrence, S. J. D., van Mourik, T., Kok, P., Koopmans, P. J., Norris, D. G., & de Lange, F. P.
755      (2018). Laminar Organization of Working Memory Signals in Human Visual Cortex.
756      *Curr Biol, 28*(21), 3435-3440 e3434. doi:10.1016/j.cub.2018.08.043
757  Liu, T., Cable, D., & Gardner, J. L. (2018). Inverted Encoding Models of Human Population
758      Response Conflate Noise and Neural Tuning Width. *J Neurosci, 38*(2), 398-408.
759      doi:10.1523/JNEUROSCI.2453-17.2017
760  Love, J., Selker, R., Marsman, M., Jamil, T., Dropmann, D., Verhagen, J., . . . Wagenmakers, E.
761      J. (2019). JASP - graphical statistical software for common statistical designs. *Journal of
762      Statistical Software, 88*.
763  Luck, S. J., & Vogel, E. K. (2013). Visual working memory capacity: from psychophysics and
764      neurobiology to individual differences. *Trends Cogn Sci, 17*(8), 391-400.
765      doi:10.1016/j.tics.2013.06.006

766 Mendoza-Halliday, D., Torres, S., & Martinez-Trujillo, J. C. (2014). Sharp emergence of feature-
767     selective sustained activity along the dorsal visual pathway. *Nat Neurosci, 17*(9), 1255-
768     1262. doi:10.1038/nn.3785
769 Muckli, L., & Petro, L. S. (2013). Network interactions: non-geniculate input to V1. *Curr Opin*
770     *Neurobiol, 23*(2), 195-201. doi:10.1016/j.conb.2013.01.020
771 Panichello, M. F., DePasquale, B., Pillow, J. W., & Buschman, T. J. (2019). Error-correcting
772     dynamics in visual working memory. *Nat Commun, 10*(1), 3366. doi:10.1038/s41467-
773     019-11298-3
774 Pearson, J. (2019). The human imagination: the cognitive neuroscience of visual mental imagery.
775     *Nat Rev Neurosci, 20*(10), 624-634. doi:10.1038/s41583-019-0202-9
776 Pelli, D. G. (1997). The VideoToolbox software for visual psychophysics: transforming numbers
777     into movies. *Spat Vis, 10*(4), 437-442.
778 Rademaker, R. L., Chunharas, C., & Serences, J. T. (2019). Coexisting representations of
779     sensory and mnemonic information in human visual cortex. *Nat Neurosci, 22*(8), 1336-
780     1344. doi:10.1038/s41593-019-0428-x
781 Riggall, A. C., & Postle, B. R. (2012). The Relationship between Working Memory Storage and
782     Elevated Activity as Measured with Functional Magnetic Resonance Imaging. *Journal of*
783     *Neuroscience, 32*(38), 12990-12998. doi:10.1523/Jneurosci.1892-12.2012
784 Roelfsema, P. R., & de Lange, F. P. (2016). Early Visual Cortex as a Multiscale Cognitive
785     Blackboard. *Annu Rev Vis Sci, 2*, 131-151. doi:10.1146/annurev-vision-111815-114443
786 Samaha, J., Sprague, T. C., & Postle, B. R. (2016). Decoding and Reconstructing the Focus of
787     Spatial Attention from the Topography of Alpha-band Oscillations. *J Cogn Neurosci,*
788     *28*(8), 1090-1097. doi:10.1162/jocn_a_00955
789 Samaha, J., Switzky, M., & Postle, B. R. (2019). Confidence boosts serial dependence in
790     orientation estimation. *J Vis, 19*(4), 25. doi:10.1167/19.4.25
791 Serences, J. T. (2016). Neural mechanisms of information storage in visual short-term memory.
792     *Vision Res, 128*, 53-67. doi:10.1016/j.visres.2016.09.010
793 Serences, J. T., Ester, E. F., Vogel, E. K., & Awh, E. (2009). Stimulus-specific delay activity in
794     human primary visual cortex. *Psychol Sci, 20*(2), 207-214. doi:10.1111/j.1467-
795     9280.2009.02276.x
796 Sprague, T. C., Adam, K. C. S., Foster, J. J., Rahmati, M., Sutterer, D. W., & Vo, V. A. (2018).
797     Inverted Encoding Models Assay Population-Level Stimulus Representations, Not
798     Single-Unit Neural Tuning. *eNeuro, 5*(3). doi:10.1523/ENEURO.0098-18.2018
799 van Ede, F., Chekroud, S. R., Stokes, M. G., & Nobre, A. C. (2019). Concurrent visual and
800     motor selection during visual working memory guided action. *Nat Neurosci, 22*(3), 477-
801     483. doi:10.1038/s41593-018-0335-6
802 van Kerkoerle, T., Self, M. W., & Roelfsema, P. R. (2017). Layer-specificity in the effects of
803     attention and working memory on activity in primary visual cortex. *Nat Commun, 8*,
804     13804. doi:10.1038/ncomms13804
805 Wang, L., Mruczek, R. E., Arcaro, M. J., & Kastner, S. (2015). Probabilistic Maps of Visual
806     Topography in Human Cortex. *Cereb Cortex, 25*(10), 3911-3931.
807     doi:10.1093/cercor/bhu277
808 Watson, A. B., & Pelli, D. G. (1983). QUEST: a Bayesian adaptive psychometric method.
809     *Percept Psychophys, 33*(2), 113-120. doi:10.3758/bf03202828
810 Yu, Q., Panichello, M. F., Cai, Y., Postle, B. R., & Buschman, T. J. (2020). Delay-period activity
811     in frontal, parietal, and occipital cortex tracks noise and biases in visual working
812     memory. *PLoS Biol, 18*(9), e3000854. doi:10.1371/journal.pbio.3000854

813    Yu, Q., & Shim, W. M. (2017). Occipital, parietal, and frontal cortices selectively maintain task-
814         relevant features of multi-feature objects in visual working memory. *Neuroimage, 157*,
815         97-107. doi:10.1016/j.neuroimage.2017.05.055
816    Yu, Q., & Shim, W. M. (2019). Temporal-Order-Based Attentional Priority Modulates
817         Mnemonic Representations in Parietal and Frontal Cortices. *Cereb Cortex, 29*(7), 3182-
818         3192. doi:10.1093/cercor/bhy184
819    Yu, Q., Teng, C., & Postle, B. R. (2020). Different states of priority recruit different neural
820         representations in visual working memory. *PLoS Biol, 18*(6), e3000769.
821         doi:10.1371/journal.pbio.3000769
822

**Figures and legends**

**Figure 1. Trial sequence of the fMRI task**

Participants performed a 1-item delayed recall task on oriented gratings. On different trials, they viewed a high-contrast (60%) grating embedded in noise, a low-contrast (10%) grating embedded in noise, or a null-contrast (0%) grating embedded in noise (i.e., pure noise patch). Participants were told orientation information was always presented, despite that the visibility of the gratings might differ. After a prolonged delay (8.5 s for two participants and 9.5 s for fifteen participants), they recalled the remembered orientation on an orientation wheel. Feedback was provided at the end of every response.

834
**Figure 2. Raw response distribution**

The raw response distribution of *high*, *low*, and *null* conditions, indicated by the gray histograms.

The black lines indicate the envelope of sample distribution.

**Figure 3. Time course of BOLD activity in early visual cortex and IPS**

**A**. Trial-averaged BOLD activity in the *Early Visual Cortex* ROI. **B**. Time course of BOLD activity in the *IPS* ROI. Dark blue, light blue, and gray lines correspond to the *high*, *low*, and *null* conditions, respectively. Data from the two subjects with 8.5-sec delay periods were included in the averaged results, but event labels below the x-axis represent the trial sequence for subjects with 9.5-sec delay periods only for illustration purposes in this and subsequent figures, with S, D, R representing the Sample, Delay, and Response periods, respectively. Shaded areas indicate ± 1 SEM.

849

**Figure 4. Orientation reconstructions in early visual cortex and IPS in the *high*, *low*, and**

***null* conditions.**

**A**. Time course of orientation reconstructions (indexed by the slope of reconstruction) in early

visual cortex, from left to right: in *high*, *low*, and *null* condition. **B**. Orientation reconstruction in

the late delay period (8-10 s after trial onset) in early visual cortex. **C**. Time course of orientation

reconstructions (indexed by the slope of reconstruction) in IPS, from left to right: in *high*, *low*,

and *null* condition. **D**. Orientation reconstruction in the late delay period in IPS. Black and red

lines correspond to sample and response, respectively. Black and red asterisks at the top of each

plot denote significance of reconstructions, for sample and response, respectively. Green

asterisks denote significance of difference between response and sample reconstructions. Event

labels below the x-axis represent the Sample (S), Delay (D), and Response (R) periods,

respectively. All *p*-values were corrected with False Discovery Rate (FDR) across conditions and
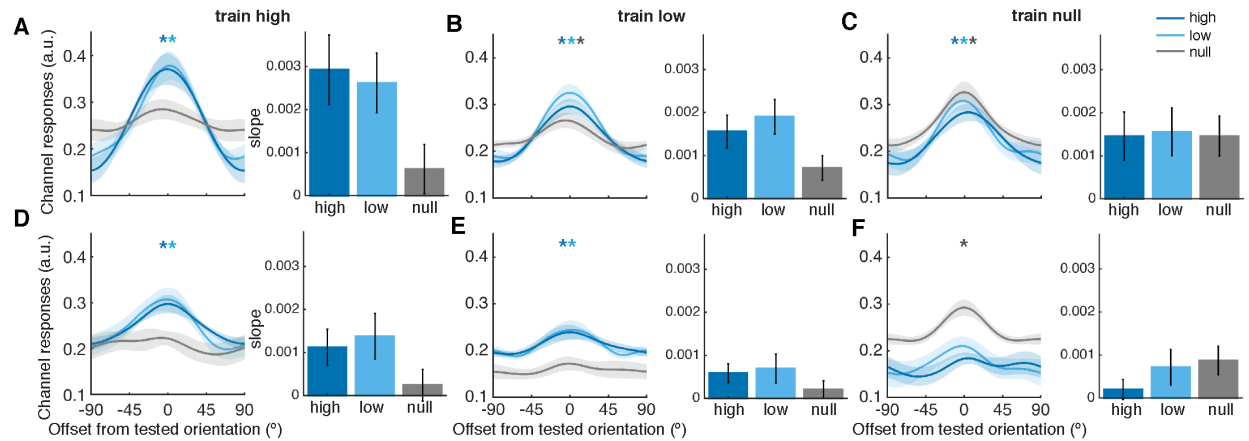
time points. Shaded areas denote ± 1 SEM. *$p < 0.05$.

864



865

**Figure 5. Time course of generalization of neural codes between the *high*, *low* and *null***

**conditions**

**A**. Time course of orientation reconstructions for *high*, *low*, and *null* conditions in early visual

cortex, from left to right: results from *high*-, *low*-, and *null*-trained IEMs. **B**. Time course of

orientation reconstructions for *high*, *low*, and *null* conditions in IPS, from left to right: results

from *high*-, *low*-, and *null*-trained IEMs. Dark blue, light blue, and gray lines correspond to the

*high*, *low*, and *null* conditions, respectively. Dark blue, light blue, and gray asterisks at the top of

each plot denote significance of reconstruction at each time point relative to baseline, for *high*,

*low*, and *null* conditions, respectively. Event labels below the x-axis represent the Sample (S),

Delay (D), and Response (R) periods, respectively. All *p*-values were corrected with False

Discovery Rate (FDR) across conditions and time points. Shaded areas indicate ± 1 SEM. **\****p* <

0.05.

878

**Figure 6. Generalization of neural codes between the *high*, *low* and *null* conditions**

**A**. Left panel: orientation reconstructions for *high*, *low*, and *null* conditions from the *high*-trained IEM, in the late delay period (8-10 s after trial onset) in early visual cortex. Right panel: slopes of these reconstructions. **B**. Orientation reconstructions and slopes for *high*, *low*, and *null* conditions from the *low*-trained IEM, in the late delay period (8-10 s after trial onset) in early visual cortex. **C**. Orientation reconstructions and slopes for *high*, *low*, and *null* conditions from the *null*-trained IEM, in the late delay period (8-10 s after trial onset) in early visual cortex. **D**. Orientation reconstructions and slopes for *high*, *low*, and *null* conditions from the *high*-trained IEM, in the late delay period (8-10 s after trial onset) in IPS. **E**. Orientation reconstructions and slopes for *high*, *low*, and *null* conditions from the *low*-trained IEM, in the late delay period (8-10 s after trial onset) in IPS. **F**. Orientation reconstructions and slopes for *high*, *low*, and *null* conditions from the *null*-trained IEM, in the late delay period (8-10 s after trial onset) in IPS. Dark blue, light blue, and gray colors correspond to the *high*, *low*, and *null* conditions, respectively. Dark blue, light blue, and gray asterisks at the top of each plot denote significance of reconstruction, for *high*, *low*, and *null* conditions, respectively. All $p$-values were corrected with False Discovery Rate (FDR) across conditions. Shaded areas and error bars indicate ± 1 SEM. *$p < 0.05$.
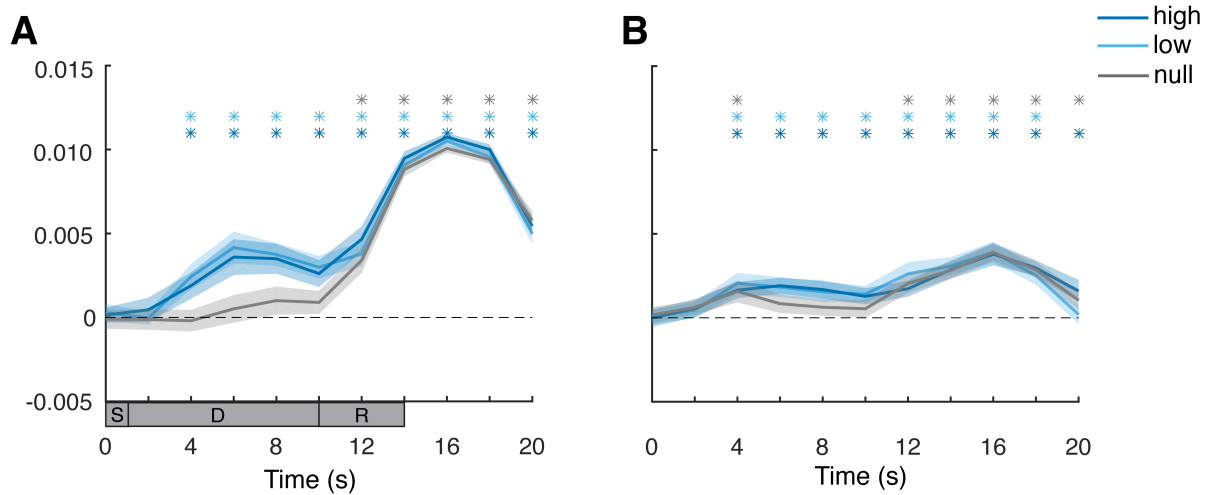
898

**Figure 7. Time course of generalization of neural codes between the *high*, *low* and *null***

**conditions, using a mixed IEM**

**A**. Time course, in early visual cortex, of orientation reconstructions for *high*, *low*, and *null*

conditions using a mixed IEM. **B**. Time course, in IPS, of orientation reconstructions for *high*,

*low*, and *null* conditions using a mixed IEM. Dark blue, light blue, and gray lines correspond to

the *high*, *low*, and *null* conditions, respectively. Dark blue, light blue, and gray asterisks at the

top of each plot denote significance of reconstruction at each time point relative to baseline, for

*high*, *low*, and *null* conditions, respectively. Event labels below the x-axis represent the Sample

(S), Delay (D), and Response (R) periods, respectively. All *p*-values were corrected with False

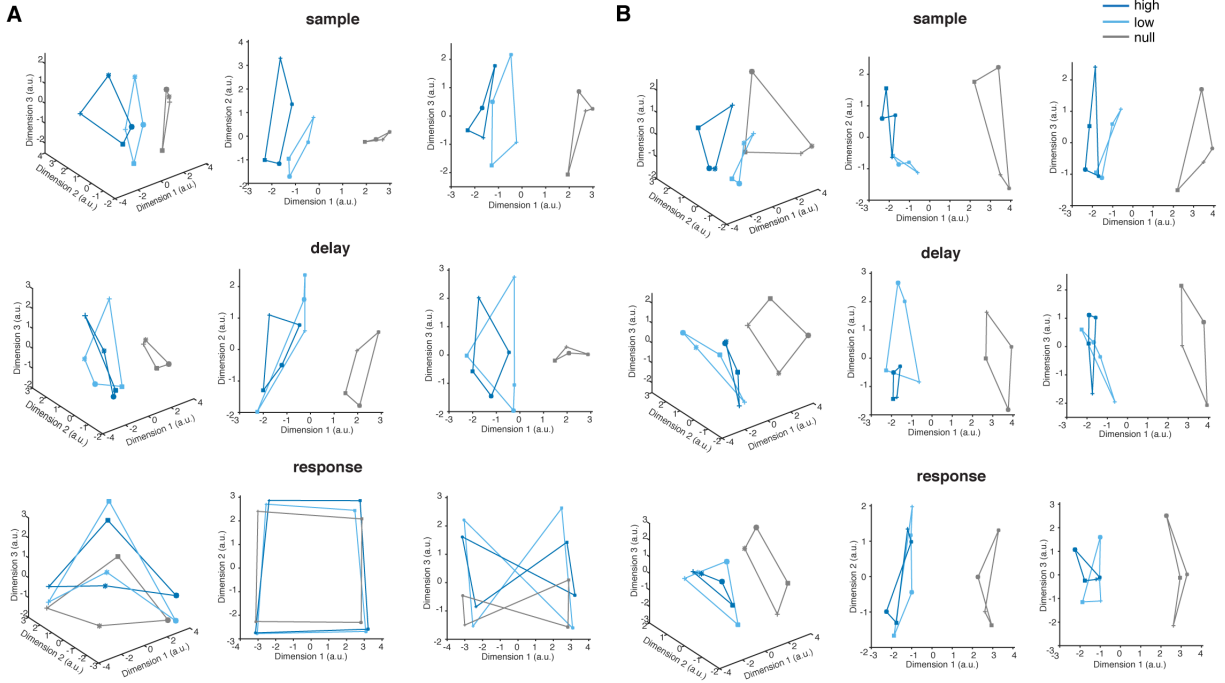Discovery Rate (FDR) across conditions and time points. Shaded areas indicate $\pm$ 1 SEM. *$p <$

0.05.

910

**Figure 8. Visualization of distances between conditions in multidimensional**

**representational space**

**A.** Visualization of representational distances between *high*, *low*, and *null* conditions, in early

visual cortex. The top panel shows the same MDS projection of data from the sample epoch from

three perspectives: 3D; a 2D view of dimension 1 vs. dimension 2; and a 2D view of dimension 1

vs. dimension. The middle panel shows three comparable views of the MDS projection of data

from the late delay period, and the bottom panel three comparable views of the MDS projection

of data from the response epoch. **B**.  MDS analyses of data from IPS, using the same display

conventions as A. Each marker represents one of the orientation bins (0-45º, 45-90º, 90-135º,

135-180º). Dark blue, light blue, and gray colors correspond to the *high*, *low*, and *null* conditions,
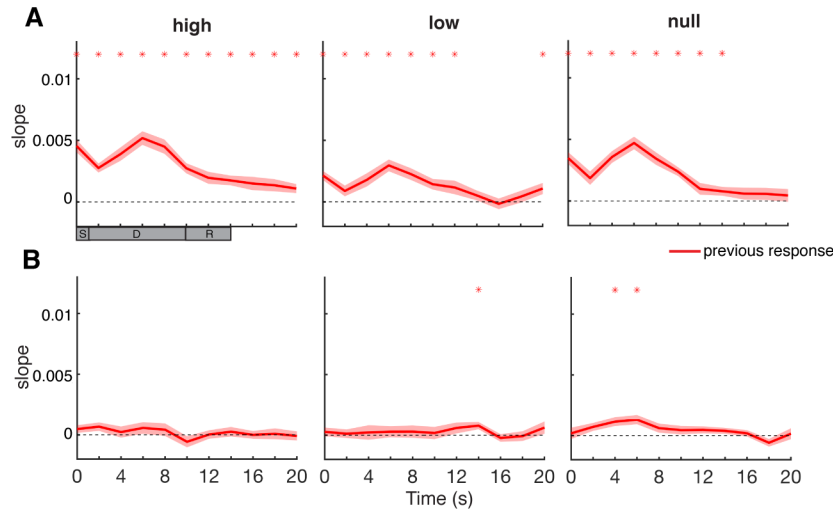
respectively. a.u.: arbitrary units.

922

**Figure 9. Reconstructions of previous-trial response in early visual cortex and IPS in the *high*, *low*, and *null* conditions.**

**A**. Time course of reconstructions of previous-trial response (indexed by the slope of reconstruction) in early visual cortex; **B**. Time course of reconstructions of previous-trial response in IPS, from left to right: in *high*, *low*, and *null* condition. Red asterisks at the top of each plot denote significance of response reconstructions. Event labels below the x-axis represent the Sample (S), Delay (D), and Response (R) periods, respectively. All $p$-values were corrected with False Discovery Rate (FDR) across conditions and time points. Shaded areas denote $\pm$ 1 SEM. *$p < 0.05$.
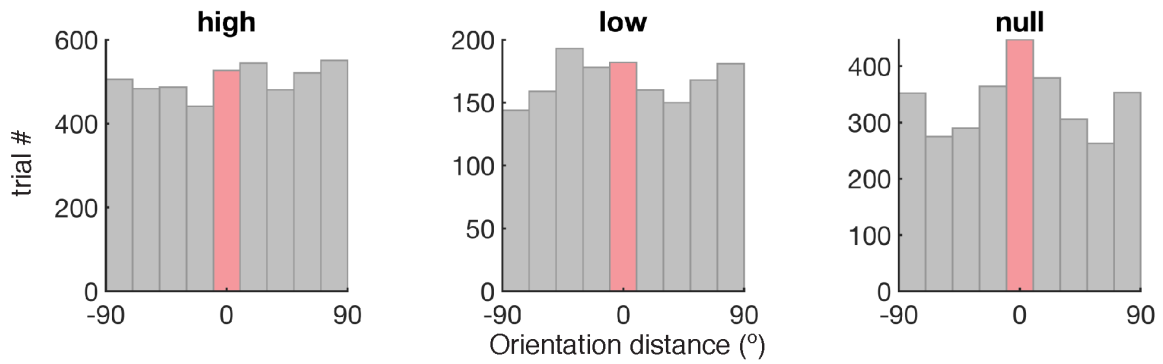
932

**Figure 10. Distribution of the distance between response on the current trial and that on the previous trial**

The distribution of the distance between response on the current trial and that on the previous trial, in the *high*, *low*, and *null* conditions, indicated by the gray histograms (bin size = 20º). The red bar at the center indicates the bin with the smallest distance (<= 10º).
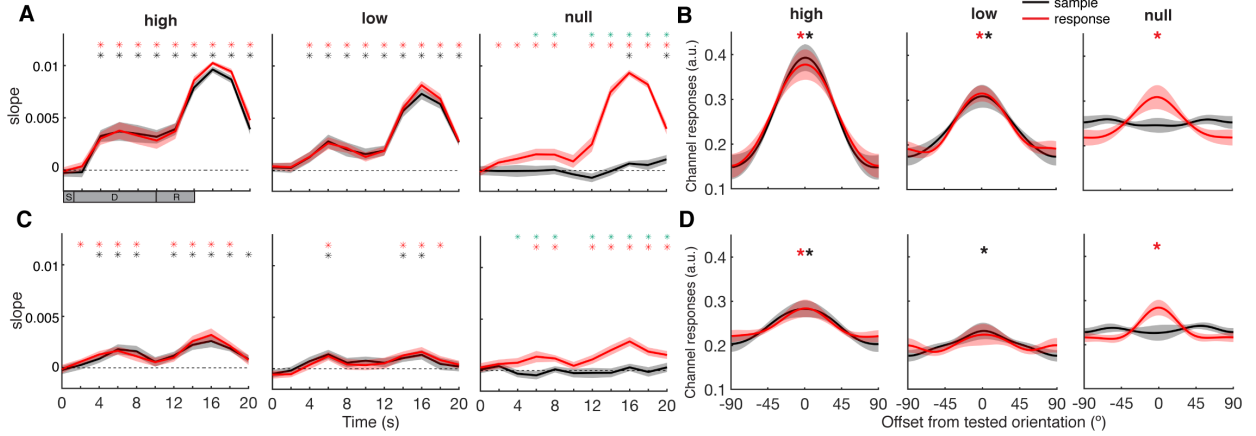
**Figure 11. Orientation reconstructions in early visual cortex and IPS in the *high*, *low*,**

**and *null* conditions, after excluding highlighted trials in Figure 10.**

**A**. Time course of orientation reconstructions (indexed by the slope of reconstruction) in

early visual cortex, after excluding trials from the bin with the smallest distance in Figure 10,

from left to right: in *high*, *low*, and *null* condition. **B**. Orientation reconstruction in the late

delay period (8-10 s after trial onset) in early visual cortex. **C**. Time course of orientation

reconstructions (indexed by the slope of reconstruction) in IPS, from left to right: in *high*,

*low*, and *null* condition. **D**. Orientation reconstruction in the late delay period in IPS. Black

and red lines correspond to sample and response, respectively. Black and red asterisks at the

top of each plot denote significance of reconstructions, for sample and response, respectively.

Green asterisks denote significance of difference between response and sample

reconstructions. Event labels below the x-axis represent the Sample (S), Delay (D), and

Response (R) periods, respectively. All *p*-values were corrected with False Discovery Rate

(FDR) across conditions and time points. Shaded areas denote ± 1 SEM. \**p* < 0.05.