1 **Title:**

2 **Spatiotemporal dynamics of self-generated imagery reveal a reverse cortical hierarchy from**

3 **cue-induced imagery**

4 **Authors:**

5 Yiheng Hu[1,2], Qing Yu[1,*]

6 **Affiliations:**

7 [1]Institute of Neuroscience, Center for Excellence in Brain Science and Intelligence Technology,

8 Chinese Academy of Sciences, Shanghai, China

9 [2]University of Chinese Academy of Sciences, Beijing, China

10

11 [*]**Correspondence should be addressed to:**

12 Qing Yu

13 Institute of Neuroscience, Center for Excellence in Brain Science and Intelligence Technology,

14 Chinese Academy of Sciences

15 Shanghai, 200031, China

16 Email: qingyu@ion.ac.cn

17

18

19

20 **Keywords**: visual imagery, frontal cortex, early visual cortex, reverse hierarchy, self-generated

21 imagery; fMRI; EEG

22

23

**Abstract**

Visual imagery, the ability to generate visual experience in the absence of direct external stimulation, allows for the construction of rich internal experience in our mental world. Most imagery studies to date have focused on cue-induced imagery, namely the to-be-imagined contents were triggered by external cues. It has remained unclear how internal experience derives volitionally in the absence of any external cues, and whether this kind of self-generated imagery relies on an analogous cortical network as cue-induced imagery. Here, leveraging a novel self-generated imagery paradigm, we systematically examined the spatiotemporal dynamics of self-generated imagery, by having participants volitionally imagining one of the orientations from a learned pool; and of cue-induced imagery, by having participants imagining line orientations based on associative cues acquired previously. Using electroencephalography (EEG) and functional magnetic resonance imaging (fMRI), in combination with multivariate encoding and decoding approaches, our results revealed largely overlapping neural signatures of cue-induced and self-generated imagery in both EEG and fMRI; yet, these neural signatures displayed substantially differential sensitivities to the two types of imagery: self-generated imagery was supported by an enhanced involvement of anterior cortex in generating and maintaining imagined contents, as evidenced by enhanced neural representations of orientations in sustained potentials in central channels in EEG, and in posterior frontal cortex in fMRI. By contrast, cue-induced imagery was supported by enhanced neural representations of orientations in alpha-band activity in posterior channels in EEG, and in early visual cortex in fMRI. These results jointly support a reverse cortical hierarchy in generating and maintaining imagery contents in self-generated versus externally-cued imagery.

**Introduction**

Visual imagery is the ability to generate visual experience from the internal world, in the absence of direct external stimulation [1]. It remains a fundamental capability of human cognition, and is central to the understanding of how our mental world is constructed. Unlike visual perception which is primarily driven by physical external stimulation and can be measured via standardized paradigms, visual imagery by definition involves cognitive processes that are ambiguous and difficult to measure in nature. Consequently, various behavioral paradigms have been used to study visual imagery; these paradigms might fundamentally differ in the exact cognitive processes involved, yet

most of them share a cue-induced nature in common, namely the contents of imagery are induced by externally presented cues. Overall, two types of imagery tasks are most frequently used: the first type of imagery task employs semantic [2-4] or associative cues [5] to trigger retrieval of imagery contents from long-term memory. In these tasks, the to-be-imagined contents are not directly accessible on the screen, but can only be inferred from long-term memory. The other type of imagery task utilizes retrocues [6,7] or mental rotation cues [8,9] to access specific memorized contents maintained or manipulated in working memory. These cue-induced imagery tasks, albeit significantly differed in their way to cue imagery, have led to several consistent observations in visual imagery: first, imagery and perception share common neural codes in early visual cortex for simple visual features [8], in object-selective high-level visual cortex for complex visual objects [2,10]. and in alpha-band activity in electroencephalography (EEG) [4], suggesting the depictive nature of visual imagery; second, neural processing during imagery follows a reverse cortical hierarchy from that during perception, which is supported by larger spatial overlap of univariate BOLD activations between imagery and perception in higher-order frontoparietal than in occipitotemporal cortex [3]; an increased top-down signal flow in imagery compared to perception, from frontal [11,12] or parietal [13] to occipital cortex; and a reversal of object representations from high-level to low-level visual cortex [2,14,15]. These findings together indicated that visual imagery involves a distributed cortical network from low-level visual cortex to higher-level visual and frontoparietal cortex [3,16], and provided empirical support for the reverse visual hierarchy model, which proposes that, as opposed to perception which triggers a feedforward sweep of neural activations along the posterior-to-anterior cortical hierarchy, imagery is initiated by top-down signals generated in higher-level cortex that trigger a cascade of neural processing in the downstream cortical areas eventually [1,17].

However, imagery experience by definition can be generated in the absence of any external stimulation, including external cues. In this context, imagery is entirely perception- or cue-independent, and the contents of imagery are self-generated from the internal world. This self-generated imagery can be regarded as a part of self-generated cognitive processes, during which an internal experience arises from intrinsic changes within an individual, rather than extrinsic changes cued from the external environment [18]. As such, self-generated imagery is much less prone to external influences, and may better reflect "pure" internally-generated mental processes. Although there have been studies on self-generated cognitive processes related to imagery, such as recalling

84  memories, envisioning the future, and mind wandering [18,19], those studies engaged complex

85  cognitive processes wherein imagery was only a part of the processes. The neural mechanism of

86  pure self-generated imagery has remained elusive. Specifically, it remains unclear whether self-

87  generated imagery works fundamentally differently from cue-induced imagery, and whether the

88  neural principles with classic cue-induced imagery paradigms would hold with self-generated

89  imagery.

90  Given that self-generated and cue-induced imagery differ primarily in the origin of imagery

91  contents, it is plausible that when participants orient internally and determine their imagery contents

92  volitionally, the reverse cortical hierarchy might be involved differently from that during cue-

93  induced imagery. A previous study has observed increased decoding performance of imagined

94  objects, as opposed to degraded decoding performance of perceived objects, from low-level to high-

95  level visual cortex [2]. The rationale is that if one brain region serves as the neural locus that initiates

96  imagery contents signals at the top of the reverse hierarchy, this region should demonstrate better

97  decoding performance of imagined than perceived contents, compared to other downstream brain

98  regions. With this logic, we would expect to see differential representational signals between self-

99  generated and cue-induced imagery along a reverse hierarchy of imagery, due to the internal origin

100  of self-generated imagery.

101  Here we set out to address these questions by comparing the neural processes underlying self-

102  generated imagery with those of cue-induced imagery. To study self-generated imagery in well-

103  controlled settings and to reduce ambiguity in imagery contents, participants' imagery contents were

104  constrained to a pool of seven fixed line orientations throughout the experiment. In self-generated

105  imagery, participants determined their imagery content freely without any associative sensory input,

106  one at a time from the seven orientations on each trial; In cue-induced imagery, participants

107  imagined one orientation based on an externally-presented associative cue, and the associations

108  between orientations and cues were learned prior to the task. We investigated the spatiotemporal

109  dynamics of these two types of imagery in a series of two experiments, using EEG (in Experiment

110  1) and fMRI (in Experiment 2), respectively. In both experiments, neural representations of imagery

111  contents were characterized using inverted encoding models (IEMs), which have been shown to be

112  a powerful tool in unveiling population-level, feature-selective representations across visual,

113  parietal, and frontal cortex in the visual working memory literature [20-22].
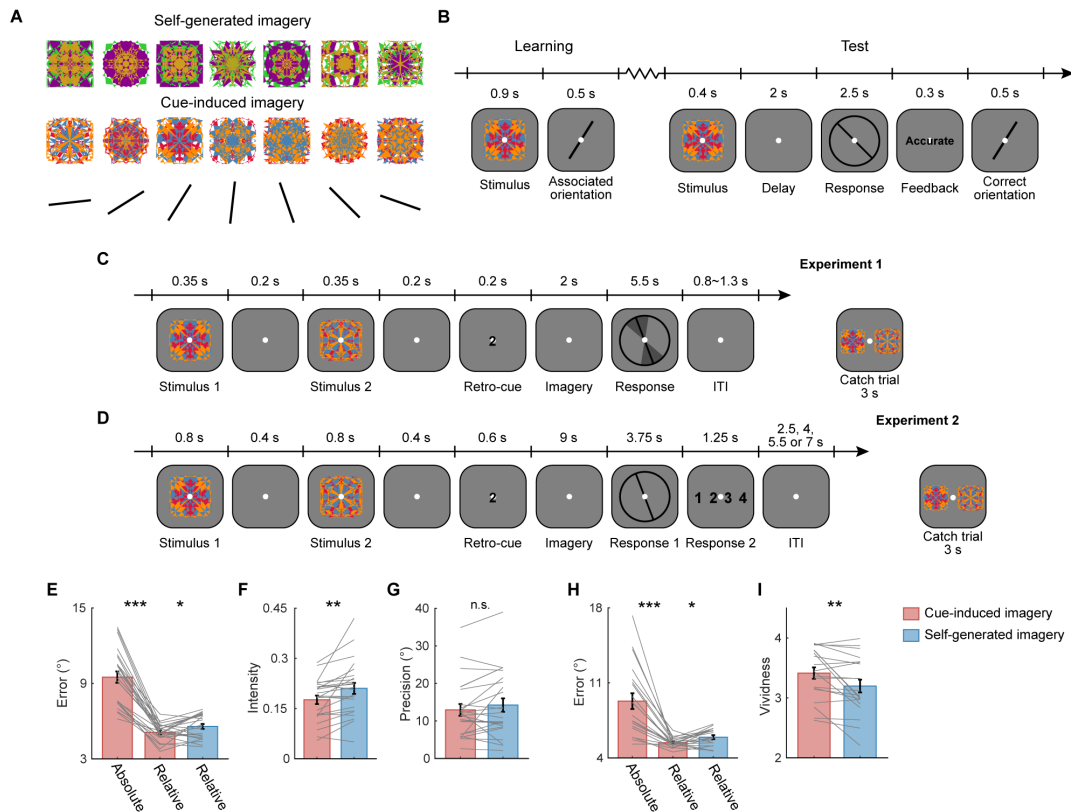
114    To preview, across two experiments, we demonstrated self-generated and cue-induced imagery

115    shared common neural representations within multiple neural signatures, while preserving

116    substantial differences in terms of the strength of representations at different levels of cortical

117    hierarchy: in EEG, enhanced orientation representations were observed in self-generated imagery

118    compared to cue-induced imagery in sustained potentials in central channels, and the opposite was

119    true in alpha-band oscillatory activity in posterior channels. In fMRI, enhanced orientation

120    representations were observed in self-generated imagery in right superior precentral sulcus (sPCS)

121    of frontal cortex, and the reverse was true in early visual cortex (EVC). In other words, the relative

122    representational strength of self-generated and cue-induced imagery also followed a frontal-to-

123    occipital reverse hierarchy. Together, these results provided the first empirical evidence, to our

124    knowledge, that frontal cortex plays a critical role in the generation and maintenance of self-

125    generated imagery contents, supporting and extending the reverse hierarchy theory of imagery.

126

127    **Results**

128    *EEG Behavior results*

129    In Experiment 1, participants performed an imagery task along with EEG recording (Figure 1A),

130    during which their imagery contents were either cued by one of seven pairs of learned associations

131    between kaleidoscope images and line orientations (Cue-induced Imagery), or self-generated from

132    the same set of seven orientations (Self-generated Imagery). During the learning session,

133    participants successfully acquired the associations between kaleidoscope images and line

134    orientations with their mean absolute recall errors being below 10°. During the EEG session,

135    participants performed the cue-induced imagery task with a mean absolute recall error of 9.47° (SD

136    = 14.28°).

**Figure 1**. Experimental paradigms and behavioral results.

A. Kaleidoscope images and line orientations used in the current study. Two sets of kaleidoscope images were used, each consisted of seven distinct images. The specific set of kaleidoscope images used for each condition (cue-induced or self-generated imagery) was counterbalanced across participants. The specific association between each kaleidoscope image and each orientation was also randomized across participants. B. Trial structure of learning and test tasks. On learning trials, participants passively viewed one kaleidoscope image, followed by its associated line orientation. On test trials, participants viewed one kaleidoscope image, and were required to report its associated orientation. Feedback was provided at the end of each trial. C-D. Trial structure of the main task. A similar trial structure was used in Experiments 1 (C) and 2 (D), and only the timing of events and the type of responses differed. Each trial began with the presentation of two consecutive kaleidoscope images followed by a retrocue. In cue-induced imagery, participants actively imagined the line orientation associated with the cued kaleidoscope image during delay; in self-generated imagery, participants freely chose one from the seven learned orientations and imagined the self-generated orientation during delay. In Experiment 1, participants reported the imagined orientation, the precision, and the intensity of their imagery; In Experiment 2, participants reported the imagined orientation and 1-4 points of vividness rating. Catch trials were interleaved to maintain participants' attention on the kaleidoscope images, and participants needed to choose the cued kaleidoscope from two probe images after retrocue. E-I. Behavioral performance in Experiments 1 and 2. E. Results of mean recall error in each condition (absolute recall error in cue-induced imagery, and relative recall error in both conditions) of Experiment 1. Colored bars indicate group mean (error bars denote ±1 SEM), gray lines indicated results from individual participants. Asterisks on top denote significance of pairwise comparisons between conditions, n.s., not significant, *: $p < 0.05$, **: $p < 0.01$, ***: $p$

161 &lt; 0.001. F. Same as E, but with results of intensity of imagery experience in Experiment 1. G. Same

162 as F, but with results of precision of imagery experience in Experiment 1. H. Same as E, but with

163 results from Experiment 2. I. Same as E, but with results of vividness rating of Experiment 2.

164

165   Because there were no correct answers on self-generated imagery trials, to compare the

166 behavioral performance between conditions, we took the least circular distance of responses to the

167 seven specific orientations as relative recall errors in both conditions. We first showed that relative

168 and absolute errors correlated with each other in cue-induced imagery ($r = 0.54$, $p = 0.006$; Figure

169 S1A), suggesting relative error may be treated as an approximation of absolute error when the latter

170 was not available in self-generated imagery. Meanwhile, relative error was significantly smaller

171 than absolute error in cue-induced imagery, $t(23) = 11.25$, $p < 0.001$. When comparing relative error

172 between conditions, we found that the mean relative error in cue-induced imagery ($5.11° \pm 3.58°$)

173 was slightly but significantly smaller than that in self-generated imagery ($5.58° \pm 3.72°$), $t(23) =$

174 $2.49$, $p = 0.020$ (Figure 1E). Furthermore, because participants were required to randomly select one

175 from seven learned orientations in self-generated imagery, we examined whether participants'

176 responses were biased towards specific orientation bins. We binned all responses into seven bins,

177 each centered at one of the seven orientations (Figure S2A, S2B). We observed a slight bias in

178 participants' response distribution in both conditions. To avoid potential influence of these biases

179 on subsequent neural analyses, we balanced the number of trials within each response bin for all

180 neural analyses (see Methods for details). Lastly, we confirmed that participants did not respond by

181 simply entering the initial orientation of the response wheel (Figure S2C, S2D). Together, these

182 results suggested that participants faithfully followed task instructions and randomly selected one

183 from seven orientations in self-generated imagery.

184   Besides recall errors, participants were also measured on the vividness of their imagery, by

185 reporting both the precision (as characterized by the angle of the response wedge) and the intensity

186 (as characterized by the darkness of the response wedge) of their imagery experience. Overall,

187 participants' subjective experience was more vivid in cue-induced imagery than in self-generated

188 imagery: participants reported a more intense imagery experience in cue-induced imagery ($0.18 \pm$

189 $0.13$) compared to in self-generated imagery ($0.21 \pm 0.14$) condition, $t(23) = 3.37$, $p = 0.003$ (Figure

190 1F). On the contrary, difference in precision between conditions was numerically but not statistically

191 different ($12.96° \pm 11.47°$ in cue-induced imagery; $14.24° \pm 12.48°$ in self-generated imagery; $t(23)$

192    = 1.43, $p$ = 0.17; Figure 1G). These results indicated that self-generated imagery produced an

193    attenuation of subjective experience in terms of subjective intensity, but not in subjective precision.

194    Moreover, in self-generated imagery, precision significantly correlated with intensity ($r$ = 0.58, $p$ =

195    0.003, Figure S1C) and relative errors ($r$ = 0.40, $p$ = 0.050, Figure S1D) across participants, while

196    the correlation between intensity and relative errors was not significant ($r$ = 0.07, $p$ = 0.747, Figure

197    S1H). Follow-up stepwise regression analysis confirmed that intensity and relative error explained

198    distinct variance in precision ($p$s = 0.002 and 0.032, respectively). In comparison, no correlation

199    was observed between any two of the behavioral measures in cue-induced imagery ($r$s< 0.40, $p$s >
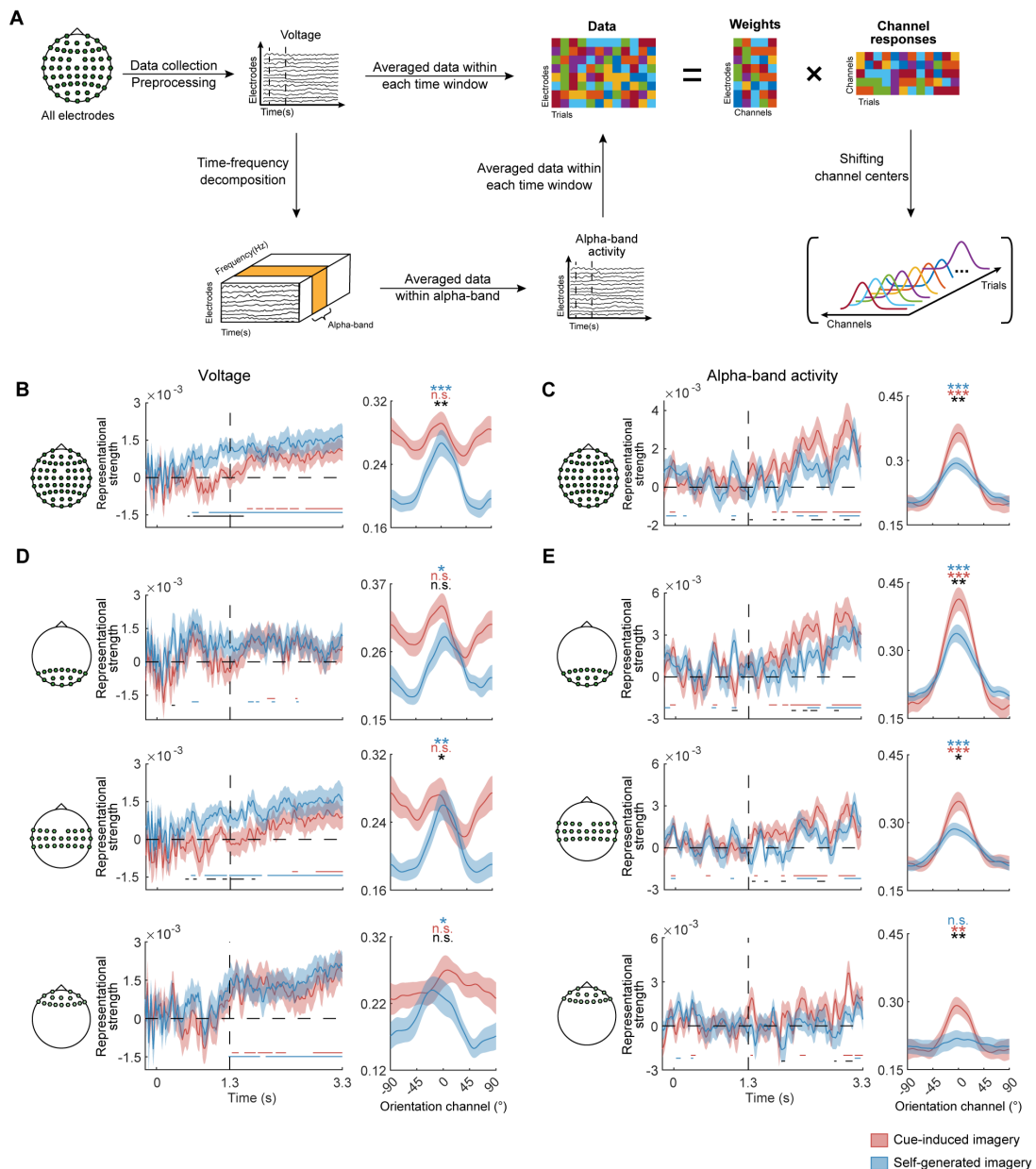
200    0.05).

201

202    ***Sustained potentials and alpha-band oscillatory activity showed differential sensitivity to self-***

203    ***generated and cue-induced imagery***

204    Having established that participants could faithfully perform the self-generated imagery task

205    following task instructions, we next seek to investigate neural signals that could potentially

206    distinguish self-generated from cue-induced imagery. For this purpose, we chose to focus on

207    stimulus-specific neural representations of imagery contents (i.e., orientations in the current study).

208    Specifically, we used participants' responses on each trial to reconstruct population-level,

209    orientation-selective representations from EEG signals using multivariate inverted encoding models

210    (IEMs). This approach has been successfully applied to investigate orientation representations in

211    various cognitive functions, for both maintenance in working memory [20,21] and retrieval from long-

212    term memory [23]. Previous studies have successfully decoded imagery contents from alpha-band [4] as

213    well as voltage signals [6] in EEG. On the other hand, recent studies on working memory [24,25] indicated

214    that alpha power and sustained potentials in EEG might reflect signals from distinct cognitive

215    processes. Given the close link between imagery and working memory [8], here we performed IEM

216    analyses on both voltage and alpha-band (8 – 12 Hz) oscillatory signals (Figure 2A), to investigate

217    whether voltage and oscillatory signals played differential roles in self-generated and cue-induced

218    imagery.

219

**Figure 2**. EEG analysis pipeline and results.

A. Pipeline of EEG analyses. Raw EEG data were collected for all electrodes. After preprocessing, preprocessed voltage data were fed into IEM analyses. Alternatively, preprocessed voltage data underwent time frequency decomposition, and the obtained power data of different frequency bands were fed into IEM analyses. B. IEM results from voltage data in all electrodes. The left panel shows time course of the strength of orientation reconstructions in cue-induced (red) and self-generated imagery (blue), from -0.2 s prior to stimulus onset until end of delay. Y axis denotes orientation representational strength, quantified using the slope of orientation reconstructions. Colored lines at the bottom denote significant time points of the corresponding condition, corrected for multiple comparisons using a cluster-based permutation method ($p < 0.01$). The vertical dashed line denotes onset of delay (at 1.3 s). The horizontal dashed line denotes baseline of reconstructions. Shaded areas denote error bars (±1 SEM). The right panel shows orientation reconstructions averaged over the selected time period of significance (0.6 – 1.7 s), in cue-induced (red) and self-generated imagery

9

234 (blue). X axis represents distance from response orientations, with 0 representing the response

235 orientation of each trial. Y axis represents reconstructed orientation channel responses in arbitrary

236 units. Colored asterisks denote significance of the corresponding condition, and black asterisk

237 denotes significance of difference between conditions. n.s., not significant, *: $p < 0.05$, **: $p < 0.01$,

238 ***: $p < 0.001$. C. same as B, but with results from alpha-band power data in all electrodes, and

239 orientation reconstructions were computed over a time window of 2-3.3 s. D. IEM results from

240 voltage data in posterior (top), central (middle), and frontal (bottom) electrodes, using the same

241 analyses and illustrations as in B. E. IEM results from alpha-band power data in posterior (top),

242 central (middle), and frontal (bottom) electrodes, using the same analyses and illustrations as in C.

243

244     Our results demonstrated that imagined orientations were represented in both voltage and

245 oscillatory signals during memory delay. Interestingly, the temporal evolution of imagery

246 representations significantly differed in these two types of signals: in voltage signals, significant

247 representations of self-generated imagined orientations ramped up around retrocue period (0.6 s

248 after trial onset) and sustained till the end of delay; whereas significant representations of cue-

249 induced imagined orientations emerged later in time and was much less stable (Figure 2B; all results

250 reported here and in subsequent analyses were corrected for multiple comparisons using a cluster-

251 based permutation method). We quantified this difference by comparing the representational

252 strength of self-generated imagery and that of cue-induced imagery during a temporal epoch around

253 retrocue (0.6 – 1.7 s after trial onset): the representational strength of orientations in self-generated

254 imagery was significant, $t(23) = 3.97$, $p = 0.0003$, and was significantly higher than that in cued-

255 induced imagery, $t(23) = 3.15$, $p = 0.002$. Meanwhile, the representational strength of orientations

256 in cue-induced imagery did not reach significance, $t(23) = 0.14$, $p = 0.447$. By contrast, in alpha-

257 band activity, significant and stable representations of both self-generated and cue-induced

258 imagined orientations ramped up around midway into the delay (2 s after trial onset; Figure 2C).

259 Moreover, when comparing the representational strength of imagery between conditions, a reversed

260 pattern was observed during late delay (2 - 3.3s): the representational strength of orientations in

261 self-generated imagery was significantly weaker than that in cue-induced imagery, $t(23) = 3.18$, $p =$

262 0.002. To validate the opposite results in voltage and oscillatory signals, we sorted participants'

263 responses into seven bins and performed multi-class classification on binned orientations using

264 support vector machines (SVMs). Representational differences between conditions in both alpha-

265 band and voltage signals remained during late delay (2 - 3.3s) with the decoding approach (Figure

266 S3). This result confirmed that the observed pattern was robust across different analytical

267     approaches used to reveal orientation representations.

268       To examine the spatial configuration of electrodes that might have contributed to the

269     representational differences between conditions, we restricted the IEM analyses to frontal, central

270     and posterior EEG electrodes, respectively. We found that, for voltage signals, only central

271     electrodes showed earlier emergence of self-generated representations, and stronger orientation

272     representations in self-generated than in cue-induced imagery ($t(23) = 2.31$, $p = 0.015$; Figure 2D),

273     suggesting that self-generated representations might primarily derive from central electrodes

274     activity. In frontal and posterior electrodes, only self-generated imagery demonstrated weak

275     orientation representations, and no difference remained in terms of either temporal dynamics or

276     representational strength between conditions (frontal: $t(23) = 1.97$, $p = 0.030$ in self-generated

277     imagery, $t(23) = 0.80$, $p = 0.217$ in difference; posterior: $t(23) = 2.12$, $p = 0.022$ in self-generated

278     imagery, $t(23) = 1.13$, $p = 0.134$ in difference). In alpha-band activity, the representational strength

279     of orientations in cue-induced imagery was higher than that in self-generated imagery in posterior

280     electrodes ($t(23) = 2.96$, $p = 0.004$; Figure 2E). Similar but weaker patterns were observed in central

281     ($t(23) = 2.21$, $p = 0.018$) and frontal electrodes ($t(23) = 3.44$, $p = 0.001$). In addition, to examine the

282     specificity of the effect to alpha-band activity, we repeated the analyses across frequencies ranging

283     from 3 to 45 Hz in posterior electrodes. We confirmed that among all frequencies, alpha-band

284     demonstrated the strongest orientation representations, as well as differences between conditions.

285     In addition, similar results were also observed in part of beta- and theta-band, but the effects were

286     overall weaker and less stable (Figure S4).

287       In Experiment 1, we demonstrated that while self-generated and cue-induced imagery shared

288     representations in both voltage and alpha-band oscillatory signals, the strength of orientation

289     representations carried in these two types of signals significantly differed between conditions in at

290     least two aspects: first, orientation representations in self-generated imagery was stronger than those

291     in cue-induced imagery in voltage signals, and the reverse was true in alpha-band signals. Second,

292     difference in voltage signals was mainly contributed by central electrodes, while difference in

293     oscillatory signals was mainly contributed by posterior electrodes. Because the trial structure of

294     these two conditions were identical, and the only difference was that imagery contents were

295     determined by different sources (self-generated versus externally-cued), we speculated that the

296     differences in signal types and spatial configurations might have reflected differences in internally-

297 generated versus externally-driven imagery: voltage signals from more anterior electrodes might

298 have reflected contents derived from self-generated imagery, and alpha oscillations from more

299 posterior electrodes might have carried information related to externally-driven processing. This

300 anterior versus posterior contrast in the spatial layout of electrodes might reflect a reverse hierarchy

301 in processing information from self-generated versus externally-cued imagery. However, due to the

302 limited spatial resolution of EEG signals, we refined our approach in a second experiment during

303 which we leveraged fMRI to investigate possible neural loci of the reverse hierarchy.

304

305 *fMRI Behavior results*

306 In order to identify brain regions that might underlie the representational differences between

307 cue-induced and self-generated imagery in voltage and alpha-band signals, we had participants

308 performed the imagery task inside an MRI scanner in Experiment 2. The procedure of Experiment

309 2 was similar to that of Experiment 1, except that the timing of events and type of responses were

310 adjusted to better suit for fMRI. Specifically, the delay period was prolonged to compensate for the

311 sluggishness of BOLD signals, and vividness rating of 1-4 points was used in order to shorten the

312 response time inside the scanner (Figure 1D).

313 The behavioral results in Experiment 2 largely replicated those in Experiment 1: the mean

314 absolute error in cue-induced imagery was 9.30° (SD = 14.15°); the mean relative errors were 5.41°

315 (SD = 3.64°) in cue-induced imagery and 5.94° (SD = 3.71°) in self-generated imagery. Relative

316 errors were significantly smaller in cue-induced compared to self-generated imagery, $t(19) = 2.78$,

317 $p = 0.012$ (Figure 1H). In terms of vividness rating, participants reported a more vivid experience

318 in cue-induced imagery (3.42 ± 0.63) than in self-generated imagery (3.20 ± 0.71), $t(19) = 2.98$, $p$

319 $= 0.008$ (Figure 1I). Moreover, vividness did not correlate with relative errors in either condition, $r$s

320 $< 0.24$, $p$s $> 0.32$. In combination with results from Experiment 1, these results together suggested

321 that precision and intensity likely reflected two different dimensions of vividness, with intensity

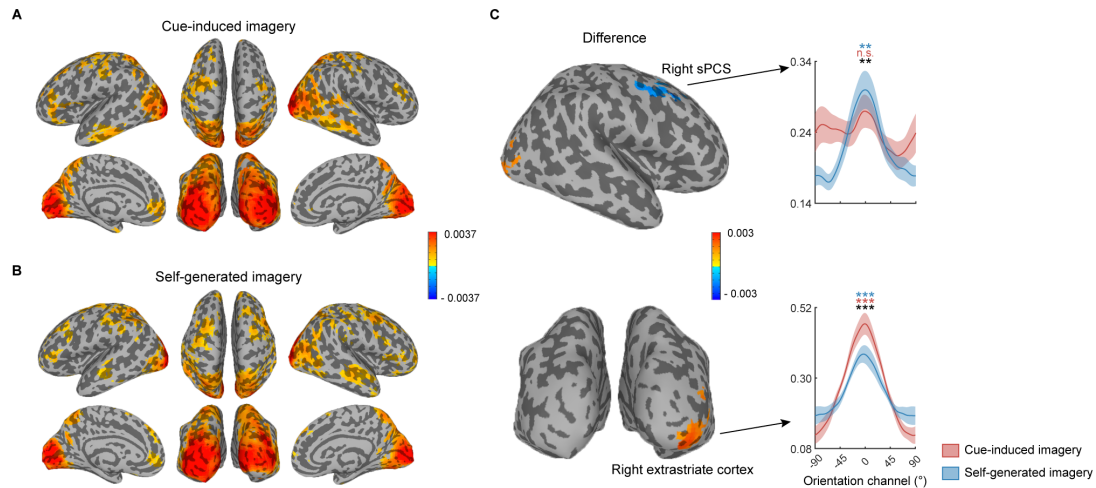322 producing qualitatively similar measures as vividness ratings.

323

324 **Whole-brain identification of representations of imagined orientations**

325 To localize brain regions showing differences in representations of imagined orientations between

326 self-generated and cue-induced imagery, we conducted a whole-brain searchlight analysis in

12

327  combination with IEM. Considering a typical hemodynamic lag of 4-6 s, the searchlight was

328  performed on data of the memory delay period (9 -12 s). Significance of searchlight results were

329  evaluated using one-tailed t-test ($p < 0.05$) and multiple comparison correction (FWE-corrected $p <$

330  0.01) to obtain the statistical parametric maps.

331  The whole-brain searchlight revealed largely overlapping brain regions for both cue-induced and

332  self-generated imagery: significant clusters with robust neural representations of imagined

333  orientations were found in a distributed network of cortical regions, including primary visual cortex

334  (V1), extrastriate cortex, intraparietal sulcus (IPS), middle and superior temporal sulcus (STS), left

335  superior precentral sulcus (sPCS), left superior medial gyrus and left middle and inferior frontal

336  gyrus. Besides these common brain regions with shared neural representations, additional clusters

337  were identified separately for the two conditions: in cue-induced imagery, imagined orientations

338  were represented in right inferior frontal sulcus (Figure 3A); in self-generated imagery, imagined

339  orientations were represented in right sPCS and right rostral lateral prefrontal cortex (rlPFC; Figure

340  3B).

341  Previous studies have revealed shared representations of perception and imagery in EVC. In the

342  current study, participants were not exposed to physical line orientations in either condition

343  throughout the imagery task. To investigate the nature of the imagery representations and to verify

344  that participants did engage visual imagery in Experiment 2, we had participants performed a

345  perception task of orientations inside the scanner following the main imagery task. We then trained

346  an IEM with perception data, and tested the perception model on imagery data in a second

347  "perception" searchlight analysis. The perception searchlight revealed similar clusters in visual

348  cortex (Figure S5), confirming a perception-like neural representation of orientations in our imagery

349  task. Additionally, similar clusters in STS for both conditions, as well as bilateral superior parietal

350  lobule and sPCS in self-generated imagery, were also identified using this perception searchlight.

**Figure 3**. Whole-brain neural representations of imagery contents in late delay.

A. Searchlight parametric map of the strength of orientation representations in late delay (9-12 s; 6 s after retrocue) in cue-induced imagery. Colors on the cortical surface denote brain regions with significant orientation representations, corrected using a cluster-based permutation method ($p < 0.01$). For demonstration purposes, clusters were thresholded at 50 voxels. B. Same as A, but with results from self-generated imagery. C. Difference map of orientation representations in A and B, with positive values denoting stronger orientation representations in cue-induced imagery, and negative values denoting stronger orientation representations in self-generated imagery. Orientation reconstructions obtained from the two significant clusters were shown, with right sPCS (top panel) demonstrating stronger orientation representations in self-generated imagery, and right extrastriate cortex demonstrating stronger orientation representations in cue-induced imagery. X axis represents distance from response orientations, with 0 representing the response orientation of each trial. Y axis represents reconstructed orientation channel responses in arbitrary units. Colored asterisks denote significance of cue-induced (red) and self-generated (blue) imagery, and black asterisk denotes significance of difference between conditions. n.s., not significant, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$.
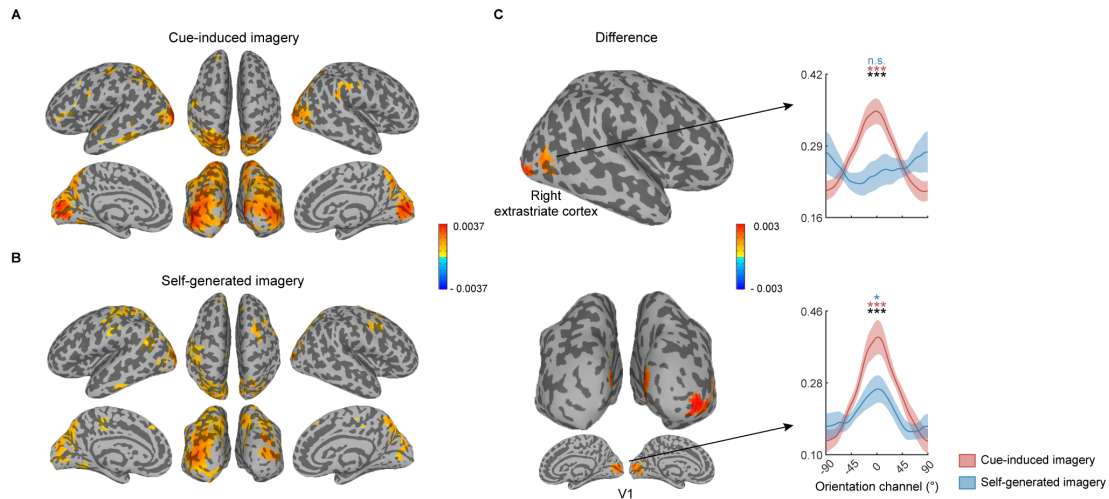
After identifying brain clusters with robust orientation representations in the two imagery conditions, we next seek to identify clusters with significant representational differences between the two. We found significant lateralization of representational differences between conditions, with right sPCS demonstrating stronger orientation representations in self-generated imagery, and right extrastriate cortex demonstrating stronger orientation representations in cue-induced imagery (Figure 3C). To better illustrate the effects, we extracted multi-voxel activation patterns from the two regions of interest (ROIs) and generated reconstructions of imagined orientations of both conditions in each ROI. The representational strength of orientations in self-generated imagery was significantly higher than that in cue-induced imagery in right sPCS ($t(19) = 3.13$, $p = 0.003$). Notably,

14

378    orientation reconstruction was significant in self-generated imagery ($t(19) = 3.54$, $p = 0.001$), but

379    not in cue-induced imagery ($t(19) = 0.69$, $p = 0.250$). In addition, the representational strength of

380    orientations in cue-induced imagery was higher than that in self-generated imagery in right

381    extrastriate cortex ($t(19) = 3.76$, $p = 0.001$), with significant orientation representations in both

382    conditions ($t(19) = 6.32$, $p < 0.001$ in cue-induced imagery, $t(19) = 3.79$, $p = 0.001$ in self-generated

383    imagery). This posterior versus anterior differences in orientation representations resembled our

384    findings in Experiment 1 which showed differential results in posterior and anterior electrodes.

385    If enhanced neural representations of orientations in sPCS supported the generation and

386    maintenance of self-generated imagery, we would anticipate the representational strength of

387    orientations in this region should be predictive of that in lower-level extrastriate cortex. Indeed,

388    Pearson correlation analysis between the two revealed significant positive correlation in self-

389    generated imagery (Figure S6B), $r = 0.48$, $p = 0.034$, but less so in cue-induced imagery (Figure

390    S6A), $r = 0.4$, $p = 0.077$.

391    Lastly, as a control, we repeated the searchlight analysis on data from an earlier epoch of the trial

392    (Figure 4, 6-9 s; 3 s after the retrocue). The representational strength of orientations in cue-induced

393    imagery was significantly higher than that in self-generated imagery in bilateral V1 ($t(19) = 3.76$, $p$

394    $= 0.001$) and right extrastriate cortex ($t(19) = 4.29$, $p < 0.001$; Figure 4C). These results were

395    consistent with a previous fMRI study demonstrating successful decoding of retrieved stimulus-

396    driven memories in early visual cortex following the onset of an associative cue [5]. On the other hand,

397    although there was no significant difference in right sPCS, there was a significant cluster in right

398    sPCS in self-generated (Figure 4B) but not in cue-induced imagery (Figure 4A), suggesting that

399    involvement of sPCS in self-generated imagery started early in the trial and became progressively

400    larger into late memory delay.

401

15

**Figure 4**. Whole-brain neural representations of imagery contents in middle delay.

A. Searchlight parametric map of the strength of orientation representations in middle delay (6-9 s; 3 s after retrocue) in cue-induced imagery. Colors on the cortical surface denote brain regions with significant orientation representations, corrected using a cluster-based permutation method ($p < 0.01$). For demonstration purposes, clusters were thresholded at 50 voxels. B. Same as A, but with results from self-generated imagery. C. Difference map of orientation representations in A and B, with positive values denoting stronger orientation representations in cue-induced imagery. Orientation reconstructions obtained from the two significant clusters were shown, both clusters showed stronger orientation representations in cue-induced imagery: in right extrastriate cortex, only orientation representations in cue-induced imagery were significant: $t(19) = 4.05$, $p < 0.001$ in cue-induced imagery, $t(19) = 0.61$, $p = 0.724$ in self-generated imagery. In V1, orientation representations in both conditions were significant, $t(19) = 4.15$, $p < 0.001$ in cue-induced imagery, $t(19) = 1.86$, $p = 0.039$ in self-generated imagery. X axis represents distance from response orientations, with 0 representing the response orientation of each trial. Y axis represents reconstructed orientation channel responses in arbitrary units. Colored asterisks denote significance of cue-induced (red) and self-generated (blue) imagery, and black asterisk denotes significance of difference between conditions. n.s., not significant, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$.

## Discussion

Even in the absence of any external stimulation, people can still self-generate contents in visual imagery. How would the neural underpinnings of self-generated imagery differ from those of classic cue-induced imagery? Here, we investigated (1) the temporal dynamics of self-generated imagery in an EEG experiment and (2) the spatial layouts of neural representations in self-generated imagery in an fMRI experiment, and contrasted the spatiotemporal dynamics of self-generated imagery with those of cue-induced imagery. Our results revealed an enhanced involvement of frontal cortex in generating and maintaining contents in self-generated imagery, as evidenced by enhanced neural

429   representations of orientations in sustained potentials in central channels in EEG, and in sPCS of

430   frontal cortex in fMRI. By contrast, cue-induced imagery was supported by enhanced neural

431   representations of orientations in alpha-band activity in posterior channels in EEG, and in visual

432   cortex in fMRI. Taken together, these results jointly support a reverse cortical hierarchy in

433   representing imagery contents in self-generated versus externally-cued imagery.

434       Previous work on visual imagery has mostly utilized cue-induced paradigms, namely the to-be-

435   imagined contents were guided by externally presented cues, either from long-term [2-5] or working

436   memory [6-9]. One advantage of using cue-induced imagery paradigms is that contents of imagery can

437   be better controlled, compared to uncontrollable situations such as mind wandering. However, real-

438   life imagery often requires imagery contents to be generated freely of external controls, yet the

439   neural mechanisms of self-generated imagery have remained largely unexplored due to limitations

440   in experimental paradigms. Although there have been several recent attempts to tackle on a related

441   problem [26,27], it has remained unclear how contents of self-generated imagery were generated and

442   how self-generated imagery differed from classic cue-induced imagery. To balance the needs for

443   well-controlled experimental settings as well as for self-generating imagery contents, here we

444   designed a novel experimental paradigm to investigate the neural mechanism of self-generated

445   imagery: participants decided freely what to imagine on each trial, with the constraint that imagery

446   contents were limited to a set of seven pre-learned line orientations. Combining the new behavioral

447   paradigm with multivariate inverted encoding models allowed us to examine the neural

448   representations of self-generated imagery contents (orientations in our case), and how these

449   representations differed from those of cue-induced imagery, beyond univariate activation changes

450   between conditions. The present study revealed several distinctive features of self-generated

451   imagery: behaviorally, the vividness of self-generated imagery was significantly reduced compared

452   to that of cue-induced imagery; neurally, self-generated imagery shared representational codes with

453   perception as well as with cue-induced imagery in early visual cortex; more interestingly,

454   converging evidence from EEG and fMRI suggested enhanced orientation representations in frontal

455   cortex in self-generated compared to cue-induced imagery. We interpreted these representational

456   differences as reflecting a reverse cortical hierarchy in representing imagery contents that were

457   generated either via internal drives or external cues. Below we discuss our findings in EEG and

458   fMRI in more details:

17

459    According to the reverse visual hierarchy model, imagery signals are initiated in the more anterior

460    part of cortex such as frontal cortex, and the signals trigger a cascade of neural processing along the

461    anterior-to-posterior cortical hierarchy [1,17]. Because the initiation signals in self-generated and cue-

462    induced imagery derived from completely different origins, we hypothesized that anterior cortex

463    would act differently in self-generated and cue-induced imagery. Our results from both EEG and

464    fMRI supported this notion: in EEG, we found that imagery contents were decodable in sustained

465    potentials in central but not in posterior electrodes in both conditions, and more importantly,

466    orientation representations emerged earlier in time in self-generated imagery, and remained stronger

467    than those in cue-induced imagery in central electrodes. Due to the poor spatial resolution of EEG

468    signals, we next turned to fMRI for the neural loci of such representational differences. Consistent

469    with the EEG findings, we observed that right sPCS of frontal cortex maintained robust orientation

470    representations of self-generated imagery but not cue-induced imagery, in both middle and late delay

471    periods. Together, these results indicated anterior cortex, especially right sPCS in the current study,

472    might serve as the critical neural locus that initiates and maintains contents in self-generated imagery.

473    The results of sPCS in the current study are broadly in line with previous work implicating a role

474    of sPCS in visual working memory [20-22]. Our results extended this finding to the imagery domain,

475    and more specifically, we demonstrated that sPCS contributed to self-generated imagery in a way

476    that was specific to imagined stimuli. Recent debates in the field of working memory have argued

477    about the specific role of higher-order frontal cortex in working memory maintenance [28,29], partially

478    due to the fact that stimulus-specific representations observed in frontal cortex during working

479    memory were substantially more variable compared to low-level visual cortex. Our finding added

480    new insights into this line of research, by demonstrating that stimulus-specific representations in

481    frontal cortex were enhanced when the level of "internality" increased as in self-generated imagery.

482    In other words, our work clearly indicated the origin of stimulus-specific representations in sPCS

483    was internal rather than external. Moreover, we demonstrated significant functional coupling

484    between stimulus-specific representations in sPCS and those in EVC in self-generated imagery, in

485    support of the view that sPCS exerts top-down control over lower-level visual cortex [30,31].

486    Intriguingly, in another recent work from our lab (unpublished), we have identified a similar reverse

487    hierarchy between EVC and IPS, for cue-induced imagery as compared to perception. The fact that

488    the top node of the reverse hierarchy moved more anteriorly from IPS to sPCS, when imagery

489 contents became more "internally-generated", possibly implied a flexible reverse hierarchy that

490 depends on the "internality" of the specific cognitive process. Last but not least, it should be noted

491 that the function of sPCS can be better understood when taking into account the type of imagined

492 stimuli used in the current study. sPCS shows robust stimulus representations for spatial or space-

493 related stimuli such as locations and orientations [20-22], but less so for non-spatial stimuli such as

494 colors [21,32]. Whether there remains a more domain-general, stimulus-nonspecific brain region in

495 self-generated imagery requires further future work to elaborate on.

496      Turning to lower-level visual cortex, imagery and perception have been shown to share neural

497 codes in early visual cortex [8,33]. Relatedly, a recent EEG experiment found imagery and perception

498 shared neural representations in alpha-band oscillatory activity in posterior electrodes [4]. In our EEG

499 experiment, although we failed to find orientation representations from sustained potentials in

500 posterior electrodes, we found imagery contents were indeed represented in alpha-band oscillatory

501 activity in posterior electrodes. Interestingly, orientation representations in alpha-band activity

502 demonstrated a reversed pattern from those in sustained potentials, and orientation representations

503 were stronger in cue-induced imagery rather than in self-generated imagery. Our fMRI searchlight

504 also demonstrated an analogous pattern that first emerged in V1 and moved onto extrastriate cortex

505 in late delay, and orientation representations of both self-generated and cue-induced imagery shared

506 common neural codes with perception in visual cortex. Moreover, the emergence of orientation

507 representations in posterior electrodes was much later in time, compared to those in central

508 electrodes in sustained potentials. Given that alpha-band oscillations carry feedback information [34-

509 36], it was likely that the orientation representations in alpha-band received feedback modulations

510 from higher-order areas, possibly frontal cortex.

511      How should we interpret the reversed patterns of results in alpha-band activity in EEG as well as

512 in visual cortex in fMRI? We noticed that this neural result echoed the behavioral difference in

513 vividness, that is, imagery experience was reported to be more vivid in cue-induced imagery than

514 in self-generated imagery. One explanation for the reduced representational strength in self-

515 generated imagery would be attenuated sensory processing for self-generated imagery, similar as

516 self-generated sensations that felt less salient than externally generated sensations [37-39], and as other

517 higher-level self-generated cognitive processes such as motor imagery, inner speech, and

518 numerosity estimation [40-42]. Attenuation in self-generated imagery can be accommodated within the

519    internal feedforward model framework [43]. This model proposes any action people take is followed

520    by a corollary discharge, which is used to predict sensory consequences of the action. When the

521    prediction matches the actual sensory feedback, the sensory consequences are attenuated. In self-

522    generated imagery, the prediction of imagery was always in line with self-generated contents, thus

523    leading to weaker neural representations in sensory cortex. Consequently, the reduction in subjective

524    experience of imagery might derive from this representational attenuation.

525        Alternatively, these results could be accounted for by the sensorimotor recruitment hypothesis,

526    which proposes that visual cortex is engaged in both perceiving external stimuli and maintaining

527    mental images [44], with shared representations between perception and working memory [45], between

528    long-term memory and perception [5], and between long-term memory and working memory [46]. In

529    these studies, significant neural representations of contents in long-term memory in early visual

530    cortex can be explained by a neural reinstatement of the to-be-retrieved information from long-term

531    memory, with the hippocampus possibly acting as the source of the modulatory signals [5]. Because

532    both of the imagery tasks in our current study engaged retrieval from long-term memory, it was

533    possible that the associative cue in cue-induced imagery facilitated memory retrieval and resulted

534    in a stronger reinstatement of imagined contents. It should be noted that the sensory attenuation and

535    sensorimotor recruitment accounts are not mutually exclusive, and might simultaneously contribute

536    to the current results.

537        It is noteworthy that in our EEG experiment, distinct result patterns were observed in sustained

538    potentials and oscillatory activity: first, orientation representations were observed in central

539    electrodes in sustained potentials, and in posterior electrodes in alpha-band activity; second, the

540    differences in representational strength between self-generated and cue-induced imagery were

541    reversed, with self-generated imagery demonstrating better orientation representation in central

542    electrodes in sustained potentials, and the opposite was true in posterior electrodes in alpha-band

543    activity in cue-induced imagery. While these results might seem difficult to explain at first glance,

544    we would like to point out that several recent studies have also reported distinct cognitive processes

545    carried by sustained potentials and oscillatory activity. For example, one study has found that

546    sustained potentials encoded contents in working memory, whereas alpha-band activity mainly

547    encoded spatial attention [24]. However, in a different memory paradigm, unattended items in working

548    memory could be decoded from alpha-band activity but not from sustained potentials [25,47]. We think

549     these results would be difficult to reconcile, without a systematic examination on the effects of the

550     two types of EEG activity along with careful experimental designs; yet, we speculated that

551     narrowing down the focus to alpha-band activity by applying wavelet transformation to sustained

552     potentials might filter out orientation-irrelevant signals and noise in other frequency bands, thereby

553     increasing the signal-noise ratio (SNR) of orientation representations in posterior electrodes. By

554     contrast, imagery-relevant signals in central electrodes might rely primarily on slow-wave cortical

555     dynamics [48], because the pattern of enhanced self-generated representations was not observed in

556     any single frequency band of oscillatory activity in central electrodes. It would be interesting for

557     future studies to examine whether the observed functional differences of frontocentral sustained

558     potentials and posterior alpha-band activity would generalize to other cognitive processes.

559     We have discussed several distinct features of self-generated imagery by comparing it with cue-

560     induced imagery; yet, there remain a few other interesting observations from the current study that

561     require further work to look into. For instance, we noticed that the majority of differences between

562     cue-induced and self-generated imagery was cortically right lateralized. Moreover, in terms of the

563     involvement of frontal cortex, there was a hint that self-generated imagery was right lateralized, and

564     cue-induced imagery was left lateralized. Whether this differential patterns in cortical lateralization

565     speaks to difference between different types of imagery remains to be further investigated [49]. In

566     addition, although we removed any potential response bias from the model training stage to avoid

567     overfitting, it would be interesting to investigate whether the two types of imagery are influenced

568     by different cognitive factors and thereby resulting in differential response bias patterns, such as the

569     oblique and attractor biases typically observed in working memory [50,51].

570     In conclusion, using both EEG and fMRI, we revealed distinctive spatiotemporal neural dynamics

571     underlying the neural basis of self-generated imagery: compared to cue-induced imagery, self-

572     generated imagery was supported by an enhanced involvement of frontal cortex, as indexed by better

573     imagery representations in sustained potentials in central channels of EEG and in sPCS of frontal

574     cortex in fMRI. The enhancement in frontal representations was accompanied by a decrease in

575     orientation representations in visual cortex, which might reflect the attenuated subjective experience

576     in vividness at the behavioral level. Research on self-generated imagery may have abundant

577     potential uses. People who suffer from schizophrenia might either have delusions that have no basis

578     in reality or generate hallucinations whose contents do not actually exist. Our results provide new

579 insights into the neural mechanisms of visual imagery, and may open up a new avenue for both

580 experimental and clinical research on imagery.

581

582 **Methods**

583 **Participants**

584 A total of forty-nine volunteers participated in the study, two of whom participated in both

585 experiments. All participants had normal or corrected-to-normal vision, reported having no

586 psychiatric or neurological disorders, provided written informed consent, and reported normal visual

587 imagery ability assessed by the Vividness of Visual Imagery Questionnaire (VVIQ) [52]. All

588 participants were recruited at Shanghai Institutes for Biological Sciences, Chinese Academy of

589 Sciences, and were monetarily compensated for their participation. The study was approved by the

590 ethical committee of Center for Excellence in Brain Science and Intelligence Technology, Chinese

591 Academy of Sciences (CEBSIT-2020028).

592 Twenty-eight volunteers participated in Experiment 1 (EEG experiment). Four participants were

593 excluded: two participants had insufficient data due to technical issues, one participant failed to

594 follow instructions, and one participant dropped out from the experiment, leaving twenty-four

595 participants in the final sample for Experiment 1 (13 females, 11 males; mean age = 24.1, SD = 2.3).

596 Twenty-three volunteers took part in Experiment 2 (fMRI experiment), all were eligible for MRI

597 scans. Three participants were excluded: two participants had insufficient data due to technical

598 issues, one participant failed to follow instructions, leaving twenty participants in the final sample

599 for Experiment 2 (11 females, 9 males; mean age = 23.6, SD = 2.3). We did not estimate sample

600 sizes for Experiment 1 or 2 a priori, but the sample size used in both experiments were comparable

601 to those in previous studies with similar approaches.

602

603 **Stimuli & Apparatus**

604 Two sets of non-semantic kaleidoscope images were used, each consisted of seven images. All of

605 them were generated by Python2 codes used in a previous study [53]. Each kaleidoscope was created

606 by overlaying three transformed hexagons. Each hexagon had a unique color and was transformed

607 by four rounds of side deflection at a random direction. The RGB values of colors were [(220,20,60),

608 (70,130,180), (255,140,0)] in one set and [(50,205,50), (139,0,139), (205,155,29)] in the other set.

609    In Experiment 1, stimuli were presented with MATLAB (R2018b, The MathWorks) and

610    Psychtoolbox 3 extensions [54,55]. They were displayed on a 48×27 cm HIKVISION LCD screen with

611    a 60 Hz refresh rate and a 1920 × 1080 resolution. The viewing distance was 62 cm. While

612    performing the task, participants' head position was stabilized by a chin rest. Responses were

613    recorded with a keyboard and a mouse.

614    In Experiment 2, All stimuli were presented using MATLAB (R2012b, The MathWorks) and

615    Psychtoolbox 3 extensions on an SINORAD LCD projector (1280 × 1024 resolution; 60 Hz refresh

616    rate). Participants viewed stimuli through a coil-mounted mirror in the scanner at a viewing distance

617    of 90.5 cm. Responses were made via two SINORAD two-key button boxes.

618

619    **Experimental paradigm and procedure**

620    **Overview**

621    The purpose of the current study was to unveil the spatiotemporal neural dynamics of self-

622    generated imagery, by contrasting with those of cue-induced imagery. In both conditions, we

623    presented kaleidoscope images instead of the actual to-be-imagined stimuli, in order to minimize

624    the influence of stimulus-driven activity in neural signals. In cue-induced imagery, imagery content

625    was determined by an external cue. Seven kaleidoscope images were used, each associated with a

626    specific line orientation. Participants were required to imagine a line with the orientation indicated

627    by the kaleidoscope image. To further eliminate stimulus-driven activity from the kaleidoscope

628    images, we adopted a retrocue imagery paradigm. On each trial, participants were presented with

629    two kaleidoscope images followed by a retrocue. The retrocue indicated the specific kaleidoscope

630    image with which the associated orientation should be imagined. In self-generated imagery, the

631    imagery content was determined by participants volitionally. Participants needed to generate their

632    imagery content on their own by freely choosing one from the seven learned orientations on each

633    trial. To match the trial time course of the cue-induced imagery, seven different kaleidoscope images

634    and a retrocue design were also used, but the kaleidoscopes were not associated with orientations.

635    The experimental paradigm was depicted in Figure 1. The specific set of kaleidoscope images used

636    for each condition was counterbalanced across participants. The specific association between each

637    kaleidoscope image and each orientation was also randomized across participants.

638    Another goal of the current study was to obtain trialwise objective and subjective measures of

639    imagery. At the end of each trial, participants were required to report their imagery content

640    (orientation) and subjective vividness. In Experiment 2, due to time limitations, we used 1-4 point

641    of vividness rating for assessing subjective vividness as used in previous studies [56]. However, the

642    standard measurement of vividness such as 1-4 point rating conflated lots of different factors of

643    subjective experience, such as subjective specificity and subjective intensity [57]. As such, to uncover

644    specific subjective experience in different dimensions, in Experiment 1 the vividness was

645    decomposed into two different dimensions: precision and intensity. Precision represented the

646    confidence in the precision of orientation report, whereas intensity indicated the subjective strength

647    of imagery content.

648

649    **Behavioral learning session**

650    Prior to the main experimental session, participants first learned the associations between seven

651    kaleidoscope images and seven specific orientations (spanning the entire orientation space and were

652    equally distant: 15°, 40.71°, 66.43°, 92.14°, 117.86°, 143.57°, 169.29°). On each trial, participants

653    passively viewed one kaleidoscope image followed by its associated line orientation. The

654    kaleidoscope image was presented for 0.9 s and then the line was shown for 0.5 s, with an inter-

655    stimulus-interval (ISI) of 0.2 s in between, followed by an inter-trial-interval (ITI) of 1.2 s. Each

656    learning block consisted of all seven association pairs presented in a randomized order. At the end

657    of each block, participants could decide whether to perform a test on their learned associations or to

658    continue with learning. Each trial started with a fixation period of 0.5 s, and then a kaleidoscope

659    image was presented for 0.4 s. After a 2-s delay, participants were required to report the

660    corresponding orientation on an orientation wheel as precisely as possible in 2.5 s. A feedback

661    message would be presented for 0.3 s, indicating whether the response was accurate (error < 5°) or

662    inaccurate (error >= 5°). After an interval of 0.2 s, the line with the correct orientation would be

663    shown for 0.5 s to consolidate memory. ITI varied in 0.8-1.3 s. The test phase consisted of 28 trials.

664    The radius of the line in both learning and testing phases varied between 3.2-4.8° on a trial-by-trial

665    basis. Participants underwent the aforementioned procedure iteratively until the mean absolute error

666    during test fell below 10°.

667

668    **Experiment 1**

669  In Experiment 1 (the EEG experiment), each trial started with the successive presentation of two

670 kaleidoscope images at the center of the screen ($3.65° \times 3.65°$ in size), each for 0.35 s with an ISI of

671 0.2 s. After 0.2 s, a retrocue followed for 0.2 s indicating which of the two stimuli should be used

672 for imagery. If images had no associations (self-generated imagery), participants needed to freely

673 choose one from seven learned orientations and imagine a line with the chosen orientation at the

674 center of the screen; if images were associated with orientations (cue-induced imagery), participants

675 needed to imagine the line at the orientation associated with the cued kaleidoscope image. After a

676 delay period of 2 s, during which participants needed to keep actively imagining the orientation,

677 participants were required to report the orientation, precision and intensity of their imagined line on

678 an orientation wheel in 5.5 s. The orientation wheel consisted of a circle with a radius of 5°, a needle

679 crossing the fixation point with the same radius, and a bowtie-shaped wedge centered on the needle.

680 The orientation of the needle represented the orientation of the imagined line, which was adjusted

681 by changing the position of the mouse cursor. The angle of wedge represented precision, which was

682 adjusted by changing the distance of cursor to the fixation. The color of the wedge indicated intensity,

683 which was adjusted by two buttons (increase or decrease) on a keyboard. The initial values of

684 orientation, angle and color were randomly chosen and participants could move the cursor and press

685 keyboard to report three variates simultaneously. Only when both operations of mouse and keyboard

686 were finished would the trial end. ITI varied in 0.8-1.3 s. The cued kaleidoscope, cued order (first

687 versus second), and condition were fully counterbalanced across blocks.

688  In each block, there were three catch trials to keep participants' attention on kaleidoscope images

689 before delay. The catch trial had the same time course as the main task trial, except that it required

690 participants to choose the cued kaleidoscope from two probe images after retrocue in 3 s. catch trial

691 ITIs were fixed at 1.05 s. At the end of each block, participants received feedback on main task

692 performance and catch trial performance. Taking account of catch trials, there were 45 trials per

693 block. All participants needed to complete 20 blocks in Experiment 1. In total, participants

694 performed 420 trials per condition. Seven participants performed the task without catch trials.

695

696 **Experiment 2**

697 **Imagery task**

698  The procedure of Experiment 2 was similar to that of Experiment 1, except that the timing of

699     events and responses were adjusted for fMRI. Participants were shown two kaleidoscope images

700     (3.26°×3.26° in size) successively, each in a 0.8-s stimulus window with a 0.4-s ISI. Then a retrocue

701     was presented for 0.6 s. During the delay period, participants imagined a line for 9 s. During the

702     response period, participants needed to rotate the needle of orientation wheel (radius = 3.7°) to

703     match the orientation of imagined line within 3.75 s and then rated their experienced vividness on a

704     scale from 1 to 4 points in 1.25 s, where 1 represented lowest vividness and 4 represented highest

705     vividness. The ITI varied in 2.5 s, 4 s, 5.5 s and 7 s. Task performance and the number of missing

706     vividness reports were provided at the end of each block as feedback. There were 16 main task trials

707     and one catch trial per block. The response time for catch trials was 3 s, and catch trial ITIs were

708     fixed at 4.5 s. Each participant completed 14 blocks in total, resulting in 112 trials per condition.

709     Three participants performed the task without catch trials.

710     **Perception task**

711     Because no physical orientations were present throughout the imagery task, in order to obtain

712     participants' neural responses to ground-truth, sensory orientations, participants completed three

713     additional blocks of perception task following the imagery task. On each trial, an oriented line whose

714     orientation was randomly chosen from seven specific orientations flickered at the center of the

715     screen for 4.5 s at a frequency of about 1.8 Hz. The radius of the line randomly varied between 3.2-

716     4.8° on a trial-by-trial basis. The ITI varied in 3 s, 4.5 s and 6 s. Participants were instructed to fixate

717     at the white fixation point and press a corresponding button whenever the fixation point turned green.

718     Each perception block consisted of 30 trials, and participants completed 90 trials in total.

719

720     **EEG recording and preprocessing**

721     EEG data were acquired using a Brain Products ActiCHamp recording system and BrainVision

722     Recorder (Brain Products GmbH, Gilching, Germany). Scalp voltage was obtained from a broad set

723     of 59 electrodes at 1000 Hz (FCz as reference). Vertical and horizontal EOG were recorded from 2

724     electrodes located ~2 cm above and below the right eye, and from 2 electrodes ~1.5 cm lateral to

725     the external canthi, respectively. Electrode impedance was kept below 30 kΩ.

726     Preprocessing analyses were performed in MATLAB (R2021a, The MathWorks) using EEGLAB

727     Toolbox [58]. The raw EEG signals were resampled at 250 Hz. Then the data were band-pass filtered

728     between 0.01 and 45 Hz. Epochs were segmented from -0.5 s to +3.6 s relative to the onset of the

729 first stimulus. The signals were baseline corrected from -0.2 s to 0 s. The epoched data were visually

730 inspected and those containing large muscle, cardiac and respiratory artifacts (except for eye blinks)

731 or extreme voltage offsets were manually removed. Independent component analysis (ICA) was

732 then performed using EEGLAB's binica algorithm for each subject to identify and remove

733 components that were associated with eye blinks [59] and eye movements [60]. Data after ICA were

734 treated as the preprocessed voltage data. To estimate oscillatory power across time and frequencies,

735 the voltage data from each channel and trial were convolved with a family of complex Morlet

736 wavelets spanning 3–45 Hz in 1 Hz steps with wavelet cycles increasing linearly between 3 and 10

737 cycles as a function of frequency. The power was calculated as the percent change of squared

738 absolute value in the resulting complex time series relative to the baseline between −0.2 s and 0 s.

739

740 **fMRI acquisition and fMRI data preprocessing**

741 MRI data were recorded using a Siemens Tim Trio 3.0 T scanner (Erlangen, Germany) with a

742 standard 32-channel phased-array head coil at the Center for Excellence in Brain Science and

743 Intelligence Technology, Chinese Academy of Sciences. Functional images were acquired with a

744 gradient echo echoplanar pulse sequence with a multiband acceleration factor of 2 (TR/TE =

745 1500/30 ms; flip angle = 60°; matrix = 74 × 74; 46 slices; voxel size = 3 mm isotropic). T1-weighted

746 anatomic images were collected using the Magnetization Prepared Rapid Acquisition Gradient Echo

747 (MPRAGE) pulse sequence (TR/TE = 2300/2.98 ms; flip angle = 9°; matrix = 256 × 256; 192 slices;

748 voxel size = 1 mm isotropic).

749 Preprocessing of MRI data was performed using AFNI [61]. The first five volumes of each

750 functional run were removed. The EPI data were then registered to the last volume of each scan

751 session and then to the T1 volume of the same session. Six nuisance regressors were included in

752 GLMs to account for head motion artifacts in six different directions. The data were then motion

753 corrected, detrended (linear, quadratic, cubic), and z-score normalized within each run.

754

755 **Quantification and statistical analyses**

756 **Behavioral data analyses**

757 In both Experiments 1 and 2, the behavioral performance of imagery could be quantified by errors

758 of the responses relative to the ground-truth orientations. In cue-induced imagery, error was

759    calculated as the circular distance between the cued and response orientations, which we referred as

760    the absolute error. In self-generated imagery, because there were no sample orientations in self-

761    generated imagery, errors in this condition were quantified by calculating the circular distance of

762    response orientations to all the seven learned orientations, and choosing the smallest error among

763    all as the relative error. As a comparison, relative errors were also computed for cue-induced imagery.

764        For vividness measurements, in Experiment 1, precision was quantified by the angle of response

765    wedge (ranging from 2° to 180°), of which the smaller angle represented smaller uncertainty of the

766    orientation of imagined line. Intensity was quantified by the grayscale value of response wedge

767    (ranging from 0 to 0.5; 0 = black, 0.5 = background color). Smaller values represented a more

768    intensive imagery experience. Thus, the smaller the values of precision and intensity, the more vivid

769    the subjective experience. In Experiment 2, the vividness rating score represented the level of

770    vividness, the larger the vividness score, the more vivid the subjective experience.

771        For each participant, means of error, precision, intensity or vividness rating in each condition

772    were calculated. We conducted two-tailed paired t tests to test the significance of the mean difference

773    between conditions. Pearson correlations were performed between precision, intensity and error to

774    assess their potential correlational relationships.

775        We took several approaches to assess whether there existed any systematic biases in participants'

776    responses. First, we examined the uniformity of response distributions in two conditions. All

777    responses were binned into seven bins, with each of the seven learned orientations as bin centers.

778    Differences in distributions were statistically assessed using $\chi^2$ tests. Second, we examined whether

779    the initial orientation had a systematic influence on response, by calculating the circular distance

780    between the initial orientation and final response orientation. Differences between bins were

781    statistically assessed using two-tailed paired t tests.

782

783    **Inverted encoding model (IEM)**

784    ***Overview***

785        All IEM analyses were performed in MATLAB using custom codes. The inverted encoding model

786    assumed that the signals in each unit (e.g., voltage or power in each electrode in EEG, or BOLD

787    signal in each voxel in fMRI) reflected the weighted sum of a small number of hypothesized feature

788    tuning channels (i.e., neuronal populations), each tuned for a different feature (i.e., orientation in

789    the current study). In our experiments, the number of hypothesized orientation tuning channels was

790    set to five (36° apart, equally spaced). We modeled the response profile of each channel to a specific

791    orientation θ as a half sinusoid raised to the 8th power (FWHM = 0.82 rad):

792
$$R = \cos(\theta - c)^8$$

793    where c was the center of the channel. Since there were no correct targets in self-generated

794    imagery, we took response orientations (round to the nearest integer) in both conditions as θ to

795    obtain the idealized responses from basis functions, which meant the θ was possible in the 1-180°

796    orientation space.

797    The IEM analysis proceeded in two stages, encoding (training) and decoding (test). We

798    partitioned our data into independent sets of training data and test data. In the encoding stage, the

799    hypothesized channel responses ($C_1$, k × n, k: the number of channels; n: the number of trials) were

800    projected to actual measured signals in training dataset ($B_1$, m × n, m: the number of units) according

801    to an unknown weight matrix (W, m × k), which could be described a general linear model of the

802    following form:

803
$$B_1 = W C_1$$

804    The weight matrix ($\hat{W}$) was obtained via least-squares estimation as follows:

805
$$\hat{W} = B_1 C_1^T (C_1 C_1^T)^{-1}$$

806    In decoding stage, the model was inverted to transform the independent test dataset ($B_2$, m × t, t:

807    the number of trials) into estimated channel responses ($C_2$, k × t) by the obtained weight matrix:

808
$$\hat{C}_2 = (\hat{W}^T \hat{W})^{-1} \hat{W}^T B_2$$

809    Following the IEM analysis in previous studies [27,62], the channel centers were not fixed but shifted

810    from 0°, 36°, 72°, 108°, 144° to 35°, 71°, 107°, 143°, 179° in 1° step for 36 iterations. We conducted

811    the above analysis in each iteration, such that all 180 orientations from 1° to 180° served as channel

812    centers. All of estimated channel responses from all iterations were combined to create responses of

813    180 orientation channels. The result, for any given orientation, can be considered a reconstruction

814    of the model's estimate of the neural representation of that orientation. This procedure ensured that

815    our reconstructions were not biased by any specific channel centers. The reconstruction of channel

816    responses was shifted to a common center (90° on x axis).

817    To characterize the strength of reconstructions, we folded the channel responses on both sides of

818    the common center, averaged them, fitted with linear regression, and then took the resulting slope

819     of linear regression as an index of the strength of reconstructions.

820

821     ***IEM procedure with all and balanced data***

822     To reveal orientation-specific neural representations of imagery in both conditions, we used

823     participants' response on each trial as the target label, and used data combined from both conditions

824     for training and testing IEMs. This mixed IEM was supposed to provide an unbiased way of making

825     comparisons between conditions [63]. To achieve this, we used a k-fold cross-validation procedure.

826     For each participant, all data from both conditions were divided into four folds. In each iteration, all

827     but one folds served as the training data, and the left-out fold served as the testing data. The

828     procedure iterated until all folds had served as training and testing data, and results from all

829     iterations were averaged, for each condition separately.

830     However, one potential drawback with the approach was that participants' responses were often

831     unbalanced across different response bins. This imbalance in trial number between response bins

832     might result in overfitting of IEM, such that orientation reconstructions from IEM might be

833     overestimated. To avoid this, we balanced trials in a way that the trial number in each of the seven

834     bins would be made equal. To be specific, we randomly drew a certain number of trials from each

835     bin, and the number of trials drawn was determined by the bin with the smallest number of trials

836     among the seven bins. This initial step would result in matched numbers of trials across bins, but

837     not necessarily between conditions. To further balance trials between conditions, we rebalanced the

838     trials by randomly removing one trial from a certain bin of the condition with more trials; and in the

839     meantime, including one trial from the same bin of the condition with fewer trials, this step was

840     iterated until the difference of trial numbers between conditions was below two. To make full use

841     of all data, the balancing and cross-validation procedures were repeated for 50 times, and the results

842     were averaged across repetitions.

843

844     ***IEM analyses with EEG***

845     In EEG experiment, 59 electrodes were divided into 3 subsets of electrodes (frontal electrodes:

846     FP1, FP2, AFz3, AFz4, AFz7, AFz8, Fz1, Fz2, Fz3, Fz4, Fz5, Fz6, Fz7, Fz8; central electrodes: FC1,

847     FC2, FC3, FC4, FC5, FC6, FT7, FT8, Cz1, Cz2, Cz3, Cz4, Cz5, Cz6, T7, T8, CPz1, CPz2, CPz3,

848     CPz4, CPz5, CPz6, TP7, TP8; posterior electrodes: Pz1, Pz2, Pz3, Pz4, Pz5, Pz6, Pz7, Pz8, POz3,

849    POz4, POz7, POz8, Oz1, Oz2). We applied IEM to voltage signals from all electrodes, frontal

850    electrodes, central electrodes and posterior electrodes separately. After time-frequency

851    decomposition of voltage data, the obtained power signals were averaged within alpha-band (8-12

852    Hz). Similarly, we performed IEM analyses on alpha-band power data, separately in global

853    electrodes and local subsets. For both voltage and power data, after balancing trials, IEM analysis

854    was performed at each time point with a sliding window of 3 time points to obtain time-resolved

855    orientation reconstructions. For IEM analyses on all frequencies, a similar procedure was conducted

856    on power data of every single frequency ranging from 3 to 45 Hz.

857

858    *IEM analyses with fMRI (Searchlight of IEM)*

859    In Experiment 2, the IEM was combined with a roving "searchlight" procedure [20,64], which

860    allowed us to reconstruct and quantify representations of imagined orientations across the entire

861    brain. For each participant, their data in the original space were warped to the MNI template [65]. We

862    used the "sphere_searchlight" class in PyMVPA toolbox [66] to perform the searchlight analysis. We

863    defined a spherical searchlight (radius = 9 mm) centered on each voxel of the whole-brain gray

864    matter mask. Considering a typical hemodynamic response lag of 4-6 s, we extracted and averaged

865    the BOLD responses in each voxel over a time period spanning 6-9 s (middle delay) and another

866    spanning 9-12 s (late delay) following the onset of stimulus, and performed IEM searchlight within

867    each time period. IEM analysis was performed using the data with all trials and balanced trials to

868    calculated the slope maps separately. Results from whole-brain searchlight were displayed on the

869    cortical surface reconstructed with FreeSurfer [67,68] and visualized with SUMA in AFNI.

870    Because no physical lines were present during the imagery task, to compare the neural

871    representation of imagery with that of perception, we performed a cross-task generalization IEM

872    analysis, by training the IEM on perception data, and testing the IEM on imagery data. We extracted

873    and averaged the responses in each voxel over a time period spanning 4.5-7.5 s of each trial of the

874    perception task to train the IEM searchlight, and tested the model on the late delay data (9-12 s) of

875    the imagery task. We did not balance trials for this analysis because samples in our perception task

876    were already balanced.

877    Note that the total trial number in each condition was smaller in Experiment 2 than in Experiment

878    1 due to prolonged trial length in fMRI studies. To avoid false positive results from an insufficient

879 number of trials, we conducted the searchlight with trial balancing and without trial balancing (i.e.,

880 using all trials), and took the intersection of these two statistical parametric maps as the final result

881 of the searchlight analyses in Experiment 2.

882

883 **Cluster-based multiple comparisons correction**

884 We used cluster-based permutation to correct for multiple comparisons across time points (in

885 EEG) and voxels (in fMRI). The overarching principle for cluster-based permutation is depicted as

886 below: all to-be-corrected test statistics were clustered in connected sets on the basis of temporal,

887 spatial, or spatiotemporal adjacency to form contiguous clusters. Cluster-level test statistics were

888 calculated by taking the sum of the test statistics within every cluster. The test statistics were then

889 permuted, and the cluster-level test statistic of the largest cluster was taken from the permuted data.

890 This procedure was repeated 10000 times to create a null distribution of test statistics. The

891 proportion of each cluster-level test statistic in true data being smaller than the cluster-level null

892 distribution was calculated as the p-value of the corresponding cluster.

893 For analyses on voltage and alpha-band data of EEG, the IEM slopes from 24 participants in each

894 condition and their paired difference between conditions were compared against zero using paired

895 or one-sample t-test to obtain the one-tailed significance ($\alpha = 0.05$) at each time point. To obtain the

896 null distribution for multiple comparisons correction, the slopes at each time point were randomly

897 multiplied by either 1 or -1 independently, and then performed one-tailed t-test across participants.

898 This procedure was repeated 10000 times, creating a null distribution of the t statistics. The

899 contiguous t statistics clusters of true data and the null distribution underwent multiple comparisons

900 correction to threshold ($\alpha = 0.01$) significant time points. For analyses on all frequency data of EEG,

901 a similar procedure was conducted on 2D time-frequency clusters.

902 For fMRI data, t-tests for the whole-brain slope map in each condition and their difference against

903 zero were performed using the AFNI function "3dttest++". To reduce the computational load in

904 permutation, we performed a two-stage procedure [69]. The randomized sign-flip procedure described

905 above was first repeated 100 times for each participant. In a second step, the randomized samples

906 were bootstrapped from each participant and then performed one-tailed t-test ($\alpha = 0.05$) across

907 participants to obtain a map of t statistics. The second step was repeated 10000 times to create the

908 null distribution. The contiguous t statistics clusters of true slope maps and the null distribution

32

909    underwent multiple comparisons correction and the obtained *p*-values were further corrected using

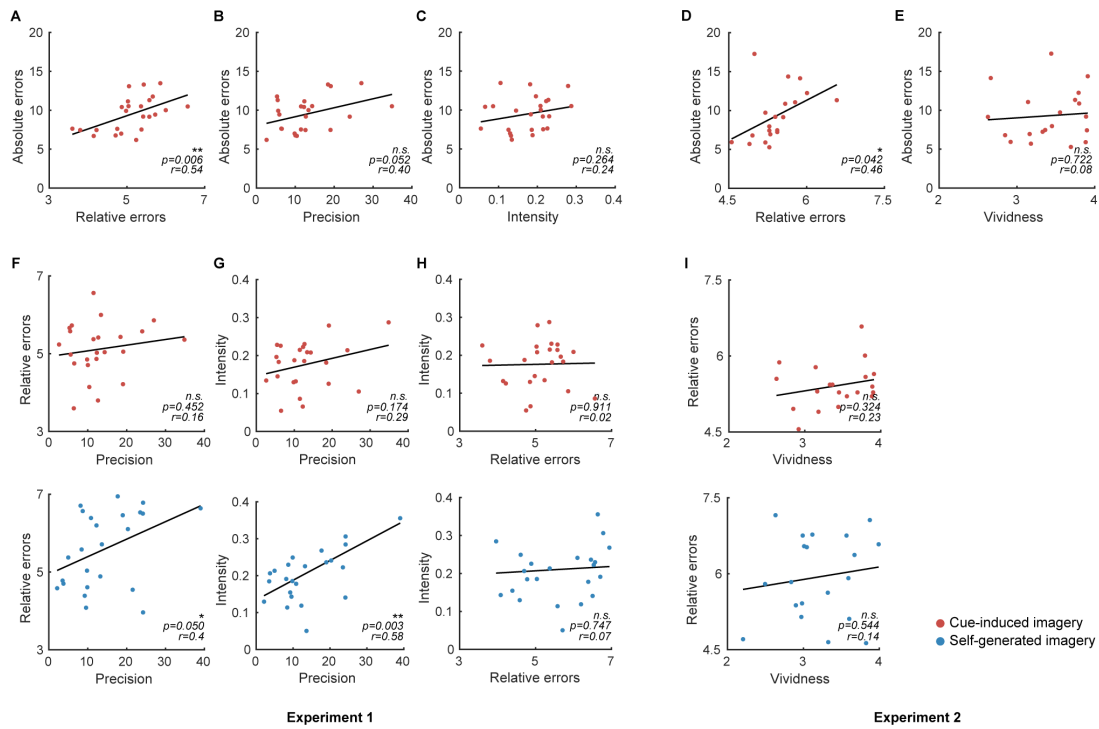910    Family-wise Error Rate (FWE) method using a threshold of α = 0.01.

911

912    **Classification using Support vector machine**

913    In order to examine whether our main results would hold with a different approach for revealing

914    orientation representations, we decoded imagery content using multi-class classification with a

915    linear support vector machine (SVM) approach. We first labeled the responses according to which

916    of the seven bins the responses belonged to. Then the trial number in each bin were balanced using

917    the method mentioned above and a k-fold (k = 4) cross-validation procedure was applied to the data

918    with balanced trials. Like the IEM analysis with balanced trials, the SVM decoding was also

919    repeated 50 times and the results from each iteration were averaged. We used the "fitcecoc()"

920    function with standard linear SVM classifier as the learner to decode the EEG data.

921

**Supplementary Figures**
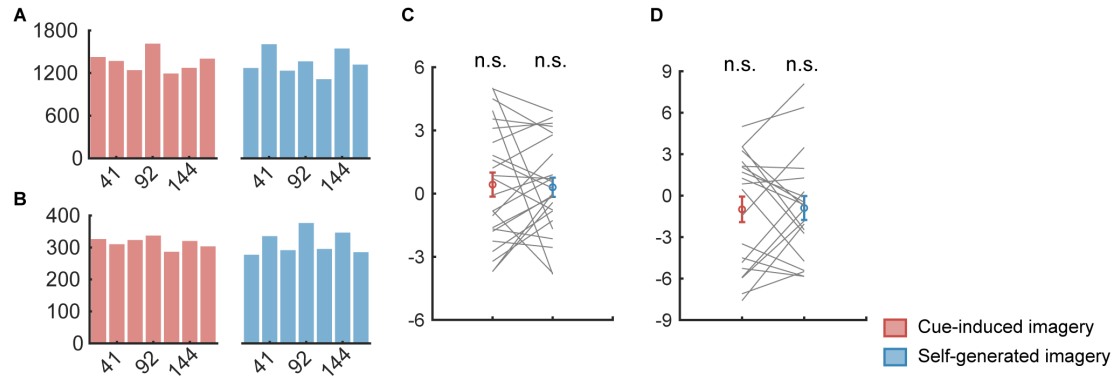


**Figure S1**. Correlations between behavioral measurements across participants.

A. Relative errors positively correlated with absolute errors in cue-induced imagery in Experiment 1. Each dot represented individual participant. Relative errors (x axis) and absolute errors (y axis) were averaged across trials for each participant. Black lines represented the best linear fit. B. Similar as A, but with correlations between precision and absolute errors in Experiment 1. C. Similar as A, but with correlations between intensity and absolute errors in Experiment 1. D. Same as A, but with results from Experiment 2. E. Similar as A, but with correlations between vividness and absolute errors in Experiment 2. F. Similar as A, but with correlations between precision and relative errors in Experiment 1. G. Similar as A, but with correlations between precision and intensity in Experiment 1. H. Similar as A, but with correlations between relative errors and intensity in Experiment 1. I. Similar as A, but with correlations between vividness and relative errors in Experiment 2. Red and blue dots represent cue-induced and self-generated imagery, respectively. Asterisks denote significance of correlations, n.s., not significant, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$.

**Figure S2**. Evaluation of response biases.

A. Histograms of response distributions, pooled from all participants in Experiment 1. x axis represents orientation bins, centered on the seven learned orientations, y axis represents frequency. Red and blue bars represent cue-induced and self-generated imagery, respectively. The uniformity of distribution was assessed using $\chi^2$ tests: cue-induced imagery: $\chi^2(23) = 86.72$, $p < 0.001$; self-generated imagery: $\chi^2(23) = 131.51$, $p < 0.001$. B. Same as a, but with results from Experiment 2: cue-induced imagery: $\chi^2(19) = 5.32$, $p = 0.503$; self-generated imagery: $\chi^2(19) = 26.26$, $p < 0.001$. C. Mean angular differences between initial probe orientations and final responses, in cue-induced (red) and self-generated (blue) imagery in Experiment 1. Colored circles indicate group mean (error bars denote $\pm 1$ SEM), gray lines indicated results from individual participants. Differences were averaged across trials for each participant and evaluated using one-sample t-test against 0: cue-induce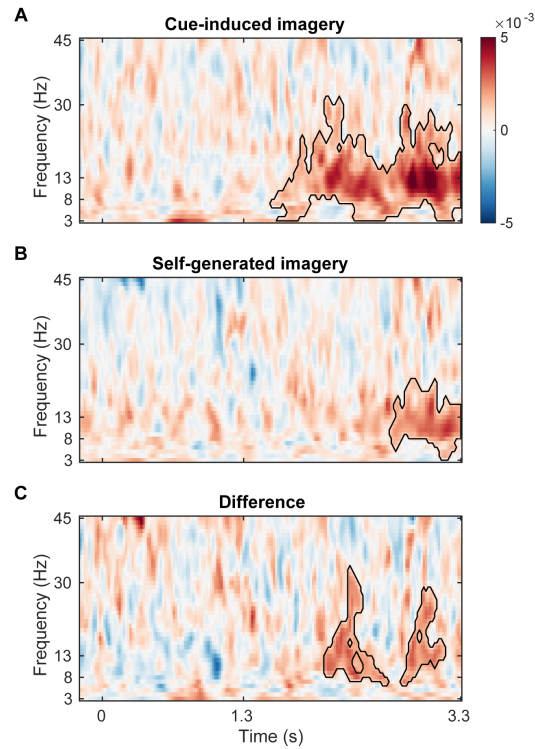d imagery: $t(23) = 0.75$, $p = 0.463$; self-generated imagery: $t(23) = 0.66$, $p = 0.518$; D. Same as C, but with results from Experiment 2: cue-induced imagery: $t(19) = 1.07$, $p = 0.297$; self-generated imagery: $t(19) = 1.04$, $p = 0.31$.
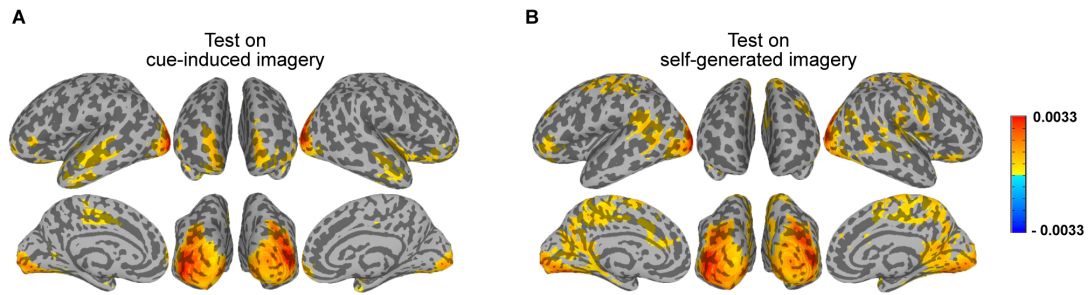
35

**Figure S3**. SVM decoding results in voltage and alpha-band activity from all electrodes. A. The left panel shows time course of decoding accuracy in cue-induced (red) and self-generated imagery (blue), from – 0.2 s prior to stimulus onset until end of delay. Y axis denotes decoding accuracy. Colored lines at the bottom denote significant time points of the corresponding condition, corrected for multiple comparisons using a cluster-based permutation method ($p < 0.01$). The vertical dashed line denotes onset of delay (at 1.3 s). The horizontal dashed line denotes chance level of 0.143 (i.e., 1/7). Shaded areas denote error bars (±1 SEM). The right panel shows decoding accuracies averaged over selected time periods of significance (2-3.3 s). Colored circles indicate group mean, error bars denote ±1 SEM, and gray lines indicated results from individual participants in cue-induced and self-generated imagery. Y axis represents decoding accuracy. Colored asterisks denote significance of the corresponding condition, and black asterisk denotes significance of difference between conditions. n.s., not significant, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$. B. same as A, but with results from alpha-band data in all electrodes.

**Figure S4**. IEM results in frequencies from 3 to 45 Hz in EEG. Orientation representational strength, as reconstructed from power data in frequencies from 3 to 45 Hz in posterior electrodes of EEG, in cue-induced imagery (top), self-generated imagery (middle), and in difference between the two (bottom). X axis denotes time, and y axis denotes frequencies. Circled areas denote significant clusters determined using a cluster-based permutation method.
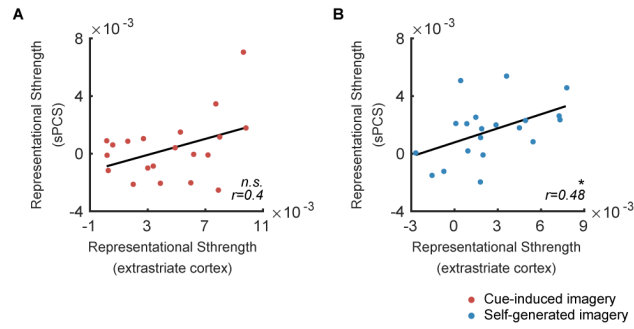
**A**

Test on
cue-induced imagery

**B**

Test on
self-generated imagery

0.0033

- 0.0033

**Figure S5**. Whole-brain neural representations of imagery contents in late delay with a perception IEM.

A. Searchlight parametric map of the strength of orientation representations in late delay (9-12 s; 6 s after retrocue) in cue-induced imagery, using an IEM trained from perception data. Colors on the cortical surface denote brain regions with significant orientation representations, corrected using a cluster-based permutation method ($p < 0.01$). For demonstration purposes, clusters were thresholded at 50 voxels. B. Same as A, but with results from self-generated imagery.

**Figure S6**. Correlations between the representational strength in right extrastriate cortex and sPCS. A. Pearson correlation between the representational strength in right extrastriate cortex and that in right sPCS in cue-induced imagery. B. Pearson correlation between the representational strength in right extrastriate cortex and that in right sPCS in self-generated imagery. Each dot represented individual participant. Representational strength in right extrastriate cortex (x axis) and representational strength in right sPCS (y axis) were averaged across trials for each participant. Black lines represented the best linear fit.

**References**

1.  Pearson, J. (2019). The human imagination: the cognitive neuroscience of visual mental imagery. Nat Rev Neurosci *20*, 624-634. 10.1038/s41583-019-0202-9.

2.  Lee, S.H., Kravitz, D.J., and Baker, C.I. (2012). Disentangling visual imagery and perception of real-world objects. Neuroimage *59*, 4064-4073. 10.1016/j.neuroimage.2011.10.055.

3.  Ganis, G., Thompson, W.L., and Kosslyn, S.M. (2004). Brain areas underlying visual mental imagery and visual perception: an fMRI study. Brain Res Cogn Brain Res *20*, 226-241. 10.1016/j.cogbrainres.2004.02.012.

4.  Xie, S., Kaiser, D., and Cichy, R.M. (2020). Visual Imagery and Perception Share Neural Representations in the Alpha Frequency Band. Curr Biol *30*, 2621-2627 e2625. 10.1016/j.cub.2020.04.074.

5.  Bosch, S.E., Jehee, J.F., Fernandez, G., and Doeller, C.F. (2014). Reinstatement of associative memories in early visual cortex is signaled by the hippocampus. J Neurosci *34*, 7493-7500. 10.1523/JNEUROSCI.0805-14.2014.

6.  Dijkstra, N., Mostert, P., Lange, F.P., Bosch, S., and van Gerven, M.A. (2018). Differential temporal dynamics during visual imagery and perception. Elife *7*. 10.7554/eLife.33904.

7.  Bainbridge, W.A., Hall, E.H., and Baker, C.I. (2021). Distinct Representational Structure and Localization for Visual Encoding and Recall during Visual Imagery. Cereb Cortex *31*, 1898-1913. 10.1093/cercor/bhaa329.

8.  Albers, A.M., Kok, P., Toni, I., Dijkerman, H.C., and de Lange, F.P. (2013). Shared representations for working memory and mental imagery in early visual cortex. Curr Biol *23*, 1427-1431. 10.1016/j.cub.2013.05.065.

9.  Iamshchinina, P., Kaiser, D., Yakupov, R., Haenelt, D., Sciarra, A., Mattern, H., Luesebrink, F., Duezel, E., Speck, O., Weiskopf, N., and Cichy, R.M. (2021). Perceived and mentally rotated contents are differentially represented in cortical depth of V1. Commun Biol *4*, 1069. 10.1038/s42003-021-02582-4.

10. Reddy, L., Tsuchiya, N., and Serre, T. (2010). Reading the mind's eye: decoding category information during mental imagery. Neuroimage *50*, 818-825. 10.1016/j.neuroimage.2009.11.084.

11. Dijkstra, N., Zeidman, P., Ondobaka, S., van Gerven, M.A.J., and Friston, K. (2017). Distinct Top-down and Bottom-up Brain Connectivity During Visual Perception and Imagery. Sci Rep *7*, 5677. 10.1038/s41598-017-05888-8.

12. Mechelli, A., Price, C.J., Friston, K.J., and Ishai, A. (2004). Where bottom-up meets top-down: neuronal interactions during perception and imagery. Cereb Cortex *14*, 1256-1265. 10.1093/cercor/bhh087.

13. Dentico, D., Cheung, B.L., Chang, J.Y., Guokas, J., Boly, M., Tononi, G., and Van Veen, B. (2014). Reversal of cortical information flow during visual imagery as compared to visual perception. Neuroimage *100*, 237-243. 10.1016/j.neuroimage.2014.05.081.

14. Dijkstra, N., Ambrogioni, L., Vidaurre, D., and van Gerven, M. (2020). Neural dynamics of perceptual inference and its reversal during imagery. Elife *9*. 10.7554/eLife.53588.

15. Horikawa, T., and Kamitani, Y. (2017). Generic decoding of seen and imagined objects using hierarchical visual features. Nat Commun *8*, 15037. 10.1038/ncomms15037.

16. Ishai, A., Ungerleider, L.G., and Haxby, J.V. (2000). Distributed neural systems for the

1039        generation of visual images. Neuron *28*, 979-990. 10.1016/s0896-6273(00)00168-9.

1040    17.    Hochstein, S., and Ahissar, M. (2002). View from the Top: Hierarchies and Reverse
1041        Hierarchies in the Visual System. Neuron *36*, 791-804. https://doi.org/10.1016/S0896-
1042        6273(02)01091-7.

1043    18.    Smallwood, J., and Schooler, J.W. (2015). The science of mind wandering: empirically
1044        navigating the stream of consciousness. Annu Rev Psychol *66*, 487-518. 10.1146/annurev-
1045        psych-010814-015331.

1046    19.    Schacter, D.L., Addis, D.R., Hassabis, D., Martin, V.C., Spreng, R.N., and Szpunar, K.K.
1047        (2012). The future of memory: remembering, imagining, and the brain. Neuron *76*, 677-
1048        694. 10.1016/j.neuron.2012.11.001.

1049    20.    Ester, E.F., Sprague, T.C., and Serences, J.T. (2015). Parietal and Frontal Cortex Encode
1050        Stimulus-Specific Mnemonic Representations during Visual Working Memory. Neuron *87*,
1051        893-905. 10.1016/j.neuron.2015.07.013.

1052    21.    Yu, Q., and Shim, W.M. (2017). Occipital, parietal, and frontal cortices selectively maintain
1053        task-relevant features of multi-feature objects in visual working memory. Neuroimage *157*,
1054        97-107. 10.1016/j.neuroimage.2017.05.055.

1055    22.    Sprague, T.C., and Serences, J.T. (2013). Attention modulates spatial priority maps in the
1056        human occipital, parietal and frontal cortices. Nat Neurosci *16*, 1879-1887.
1057        10.1038/nn.3574.

1058    23.    Sutterer, D.W., Foster, J.J., Serences, J.T., Vogel, E.K., and Awh, E. (2019). Alpha-band
1059        oscillations track the retrieval of precise spatial representations from long-term memory. J
1060        Neurophysiol *122*, 539-551. 10.1152/jn.00268.2019.

1061    24.    Bae, G.Y., and Luck, S.J. (2018). Dissociable Decoding of Spatial Attention and Working
1062        Memory from EEG Oscillations and Sustained Potentials. J Neurosci *38*, 409-422.
1063        10.1523/JNEUROSCI.2860-17.2017.

1064    25.    Barbosa, J., Lozano-Soldevilla, D., and Compte, A. (2021). Pinging the brain with visual
1065        impulses reveals electrically active, not activity-silent, working memories. PLoS Biol *19*,
1066        e3001436. 10.1371/journal.pbio.3001436

1067    10.1371/journal.pbio.3001436.g001.

1068    26.    Koenig-Robert, R., and Pearson, J. (2019). Decoding the contents and strength of imagery
1069        before volitional engagement. Scientific Reports *9*. 10.1038/s41598-019-39813-y.

1070    27.    Yu, Q., and Postle, B.R. (2021). The Neural Codes Underlying Internally Generated
1071        Representations in Visual Working Memory. J Cogn Neurosci, 1-16.
1072        10.1162/jocn_a_01702.

1073    28.    Christophel, T.B., Klink, P.C., Spitzer, B., Roelfsema, P.R., and Haynes, J.D. (2017). The
1074        Distributed Nature of Working Memory. Trends Cogn Sci *21*, 111-124.
1075        10.1016/j.tics.2016.12.007.

1076    29.    Postle, B.R., and Yu, Q. (2020). Neuroimaging and the localization of function in visual
1077        cognition. Visual Cognition *28*, 447-452. 10.1080/13506285.2020.1777237.

1078    30.    Noudoost, B., Chang, M.H., Steinmetz, N.A., and Moore, T. (2010). Top-down control of
1079        visual attention. Curr Opin Neurobiol *20*, 183-190. 10.1016/j.conb.2010.02.003.

1080    31.    Veniero, D., Gross, J., Morand, S., Duecker, F., Sack, A.T., and Thut, G. (2021). Top-down
1081        control of visual cortex by the frontal eye fields through oscillatory realignment. Nat
1082        Commun *12*, 1757. 10.1038/s41467-021-21979-7.

32. Yu, Q., and Shim, W.M. (2019). Temporal-Order-Based Attentional Priority Modulates Mnemonic Representations in Parietal and Frontal Cortices. Cereb Cortex *29*, 3182-3192. 10.1093/cercor/bhy184.

33. Ragni, F., Tucciarelli, R., Andersson, P., and Lingnau, A. (2020). Decoding stimulus identity in occipital, parietal and inferotemporal cortices during visual mental imagery. Cortex *127*, 371-387. 10.1016/j.cortex.2020.02.020.

34. van Kerkoerle, T., Self, M.W., Dagnino, B., Gariel-Mathis, M.A., Poort, J., van der Togt, C., and Roelfsema, P.R. (2014). Alpha and gamma oscillations characterize feedback and feedforward processing in monkey visual cortex. Proc Natl Acad Sci U S A *111*, 14332-14341. 10.1073/pnas.1402773111.

35. Michalareas, G., Vezoli, J., van Pelt, S., Schoffelen, J.M., Kennedy, H., and Fries, P. (2016). Alpha-Beta and Gamma Rhythms Subserve Feedback and Feedforward Influences among Human Visual Cortical Areas. Neuron *89*, 384-397. 10.1016/j.neuron.2015.12.018.

36. Bastos, A.M., Vezoli, J., Bosman, C.A., Schoffelen, J.M., Oostenveld, R., Dowdall, J.R., De Weerd, P., Kennedy, H., and Fries, P. (2015). Visual areas exert feedforward and feedback influences through distinct frequency channels. Neuron *85*, 390-401. 10.1016/j.neuron.2014.12.018.

37. Hughes, G., Desantis, A., and Waszak, F. (2013). Mechanisms of intentional binding and sensory attenuation: the role of temporal prediction, temporal control, identity prediction, and motor prediction. Psychol Bull *139*, 133-151. 10.1037/a0028566.

38. Hughes, G., and Waszak, F. (2011). ERP correlates of action effect prediction and visual sensory attenuation in voluntary action. Neuroimage *56*, 1632-1640. 10.1016/j.neuroimage.2011.02.057.

39. Benazet, M., Thenault, F., Whittingstall, K., and Bernier, P.M. (2016). Attenuation of visual reafferent signals in the parietal cortex during voluntary movement. J Neurophysiol *116*, 1831-1839. 10.1152/jn.00231.2016.

40. Jack, B.N., Le Pelley, M.E., Han, N., Harris, A.W.F., Spencer, K.M., and Whitford, T.J. (2019). Inner speech is accompanied by a temporally-precise and content-specific corollary discharge. Neuroimage *198*, 170-180. 10.1016/j.neuroimage.2019.04.038.

41. Kilteni, K., Andersson, B.J., Houborg, C., and Ehrsson, H.H. (2018). Motor imagery involves predicting the sensory consequences of the imagined movement. Nat Commun *9*, 1617. 10.1038/s41467-018-03989-0.

42. Stripeikyte, G., Pereira, M., Rognini, G., Potheegadoo, J., Blanke, O., and Faivre, N. (2021). Increased Functional Connectivity of the Intraparietal Sulcus Underlies the Attenuation of Numerosity Estimations for Self-Generated Words. J Neurosci *41*, 8917-8927. 10.1523/JNEUROSCI.3164-20.2021.

43. Miall, R.C., and Wolpert, D.M. (1996). Forward Models for Physiological Motor Control. Neural Networks *9*, 1265-1279. https://doi.org/10.1016/S0893-6080(96)00035-4.

44. D'Esposito, M., and Postle, B.R. (2015). The cognitive neuroscience of working memory. Annu Rev Psychol *66*, 115-142. 10.1146/annurev-psych-010814-015031.

45. Harrison, S.A., and Tong, F. (2009). Decoding reveals the contents of visual working memory in early visual areas. Nature *458*, 632-635. 10.1038/nature07832.

46. Vo, V.A., Sutterer, D.W., Foster, J.J., Sprague, T.C., Awh, E., and Serences, J.T. (2021). Shared Representational Formats for Information Maintained in Working Memory and

1127    Information Retrieved from Long-Term Memory. Cereb Cortex. 10.1093/cercor/bhab267.

1128  47.  Wolff, M.J., Jochim, J., Akyurek, E.G., and Stokes, M.G. (2017). Dynamic hidden states
1129        underlying working-memory-guided behavior. Nat Neurosci *20*, 864-871. 10.1038/nn.4546.

1130  48.  Hardstone, R., Flounders, M.W., Zhu, M., and He, B.J. (2022). Frequency-specific neural
1131        signatures of perceptual content and perceptual stability. eLife *11*. 10.7554/eLife.78108.

1132  49.  Liu, J., Spagna, A., and Bartolomeo, P. (2022). Hemispheric asymmetries in visual mental
1133        imagery. Brain Struct Funct *227*, 697-708. 10.1007/s00429-021-02277-w.

1134  50.  Yu, Q., Panichello, M.F., Cai, Y., Postle, B.R., and Buschman, T.J. (2020). Delay-period
1135        activity in frontal, parietal, and occipital cortex tracks noise and biases in visual working
1136        memory. PLoS Biol *18*, e3000854. 10.1371/journal.pbio.3000854.

1137  51.  Panichello, M.F., DePasquale, B., Pillow, J.W., and Buschman, T.J. (2019). Error-correcting
1138        dynamics in visual working memory. Nat Commun *10*, 3366. 10.1038/s41467-019-11298-
1139        3.

1140  52.  Marks, D.F. (1973). Visual imagery differences in the recall of pictures. Br J Psychol *64*,
1141        17-24. 10.1111/j.2044-8295.1973.tb01322.x.

1142  53.  Voss, J.L., Baym, C.L., and Paller, K.A. (2008). Accurate forced-choice recognition without
1143        awareness of memory retrieval. Learn Mem *15*, 454-459. 10.1101/lm.971208.

1144  54.  Brainard, D.H. (1997). The Psychophysics Toolbox. Spat Vis *10*, 433-436.

1145  55.  Pelli, D.G. (1997). The VideoToolbox software for visual psychophysics: transforming
1146        numbers into movies. Spat Vis *10*, 437–442.

1147  56.  Dijkstra, N., Bosch, S.E., and van Gerven, M.A. (2017). Vividness of Visual Imagery
1148        Depends on the Neural Overlap with Perception in Visual Areas. J Neurosci *37*, 1367-1373.
1149        10.1523/JNEUROSCI.3022-16.2016.

1150  57.  Fazekas, P., Nemeth, G., and Overgaard, M. (2020). Perceptual Representations and the
1151        Vividness of Stimulus-Triggered and Stimulus-Independent Experiences. Perspect Psychol
1152        Sci *15*, 1200-1213. 10.1177/1745691620924039.

1153  58.  Delorme, A., and Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of
1154        single-trial EEG dynamics including independent component analysis. J Neurosci Methods
1155        *134*, 9–21.

1156  59.  Jung, T.P., Makeig, S., Westerfield, M., Townsend, J., Courchesne, E., and Sejnowski, T.J.
1157        (2000). Removal of eye activity artifacts from visual event-related potentials in normal and
1158        clinical subjects. Clin Neurophysiol *111*, 1745-1758. 10.1016/s1388-2457(00)00386-2.

1159  60.  Drisdelle, B.L., Aubin, S., and Jolicoeur, P. (2017). Dealing with ocular artifacts on
1160        lateralized ERPs in studies of visual-spatial attention and memory: ICA correction versus
1161        epoch rejection. Psychophysiology *54*, 83-99. 10.1111/psyp.12675.

1162  61.  Cox, R.W. (1996). AFNI: software for analysis and visualization of functional magnetic
1163        resonance neuroimages. Comput Biomed Res *29*, 162-173.

1164  62.  Rademaker, R.L., Chunharas, C., and Serences, J.T. (2019). Coexisting representations of
1165        sensory and mnemonic information in human visual cortex. Nat Neurosci *22*, 1336-1344.
1166        10.1038/s41593-019-0428-x.

1167  63.  Sprague, T.C., Adam, K.C.S., Foster, J.J., Rahmati, M., Sutterer, D.W., and Vo, V.A. (2018).
1168        Inverted Encoding Models Assay Population-Level Stimulus Representations, Not Single-
1169        Unit Neural Tuning. eNeuro *5*. 10.1523/ENEURO.0098-18.2018.

1170  64.  Kriegeskorte, N., Goebel, R., and Bandettini, B. (2006). Information-based functional brain

1171          mapping. PNAS *103*, 3863-3868.

1172   65.    Fonov, V., Evans, A.C., Botteron, K., Almli, C.R., McKinstry, R.C., Collins, D.L., and

1173          Brain Development Cooperative, G. (2011). Unbiased average age-appropriate atlases for

1174          pediatric studies. Neuroimage *54*, 313-327. 10.1016/j.neuroimage.2010.07.033.

1175   66.    Hanke, M., Halchenko, Y.O., Sederberg, P.B., Hanson, S.J., Haxby, J.V., and Pollmann, S.

1176          (2009). PyMVPA: A python toolbox for multivariate pattern analysis of fMRI data.

1177          Neuroinformatics *7*, 37-53. 10.1007/s12021-008-9041-y.

1178   67.    Fischl, B., Sereno, M.I., and Dale, A.M. (1999). Cortical surface-based analysis. II:

1179          Inflation, flattening, and a surfacebased coordinate system. Neuroimage *9*, 195–207.

1180   68.    Fischl, B., Liu, A., and Dale, A.M. (2001). Automated manifold surgery: constructing

1181          geometrically accurate and topologically correct models of the human cerebral cortex.

1182          IEEE Trans Med Imaging *20*, 70-80. 10.1109/42.906426.

1183   69.    Stelzer, J., Chen, Y., and Turner, R. (2013). Statistical inference and multiple testing

1184          correction in classification-based multi-voxel pattern analysis (MVPA): random

1185          permutations and cluster size control. Neuroimage *65*, 69-82.

1186          10.1016/j.neuroimage.2012.09.063.

1187