

1 **A shared code for perceiving and imagining objects in human ventral
2 temporal cortex**

3
4 **Authors:**
5 V. S. Wadia^{1,3}, C. M. Reed², J. M. Chung², L. M. Bateman², A. N. Mamelak¹, U. Rutishauser^{ψ1,2,3,4},
6 D. Y. Tsao^{ψ5,6}
7

8 **Affiliations:**
9 1 Department of Neurosurgery, Cedars-Sinai Medical Center, Los Angeles, CA, USA
10 2 Department of Neurology, Cedars-Sinai Medical Center, Los Angeles, CA, USA
11 3 Division of Biology and Bioengineering, California Institute of Technology, Pasadena, CA, USA
12 4 Center for Neural Science and Medicine, Department of Biomedical Sciences, Cedars-Sinai
13 Medical Center, Los Angeles, CA, USA
14 5 Department of Neuroscience, Helen Wills Neuroscience Institute, University of California,
15 Berkeley, CA, USA
16 6 Howard Hughes Medical Institute, Berkeley, CA, USA
17
18 ψ - Joint senior authors
19

20 **Keywords:**
21 Visual perception; object recognition; visual imagery; inferotemporal cortex; intracranial
22 electrophysiology; generative model; human brain;
23

24 **Abstract:** Mental imagery is a remarkable phenomenon that allows us to remember previous
25 experiences and imagine new ones. Animal studies have yielded rich insight into mechanisms for
26 visual perception, but the neural mechanisms for visual imagery remain poorly understood. Here,
27 we first determined that ~80% of visually responsive single neurons in human ventral temporal
28 cortex (VTC) used a distributed axis code to represent objects. We then used that code to
29 reconstruct objects and generate maximally effective synthetic stimuli. Finally, we recorded
30 responses from the same neural population while subjects imagined specific objects and found
31 that ~40% of axis-tuned VTC neurons recapitulated the visual code. Our findings reveal that visual
32 imagery is supported by reactivation of the same neurons involved in perception, providing single
33 neuron evidence for existence of a generative model in human VTC.
34

35 **One Sentence Summary:** Single neurons in human temporal cortex use feature axes to
36 encode objects, and imagery reactivates this code.
37
38
39
40

41 **Introduction**

42

43 Mental imagery refers to our brains' capacity to generate percepts, emotions, and thoughts in
44 the absence of external stimuli. This ability pervades many aspects of the human condition. It
45 allows for the generation of visual art (1–3), musical composition (4–7), and creative writing (8–
46 10). It subserves efficient planning (11, 12) and navigation (13–15) via the simulation of actions
47 and outcomes. It is also the basis for calling to mind a recent experience, person, place, or object,
48 which is a key aspect of episodic memory (11, 16–20). In a clinical setting, uncontrolled mental
49 imagery can contribute to psychological disorders including anxiety, schizophrenia, and post-
50 traumatic stress disorder (21).

51

52 Perhaps the most consistent and established finding in mental imagery research is that imagery
53 of a given sense co-opts the neural machinery used for perception, in other words imagery is
54 supported by activation of sensory areas. This has been shown explicitly during auditory (22–25),
55 olfactory (26, 27), tactile (28), speech (29), and even motor imagery (30–34), though most
56 extensively in visual imagery (35–46); damage in or loss of either dorsal or ventral visual pathways
57 (47) often leads to parallel deficits in visual imagery (48–50). However, these studies lack the
58 spatial resolution to determine whether imagery reactivates the exact same neural populations
59 that support perception, i.e., the feedforward neural code (which would imply a generative
60 model) or whether instead imagery is subserved by activation of separate circuitry roughly
61 located in the same cortical regions. The concept of a generative model wherein representations
62 can be hierarchically re-activated in a top-down manner has great power computationally,
63 explaining self-supervised learning (51), inference under ambiguity (52), and object-based
64 attention (53). Yet, direct single-neuron evidence for regeneration of the sensory code during
65 imagery has been conspicuously lacking.

66

67 Here, we attempt to shed light on the single-neuron mechanisms of visual imagery by
68 determining the code for visual objects and then examining whether this code is reactivated
69 during imagery. We focused our investigations on ventral temporal cortex (often referred to as
70 inferotemporal cortex), a large swath of the primate temporal lobe dedicated to representing
71 visual objects (54, 55). We recorded single neurons in human patients implanted to localize their
72 focal epilepsy (56) as they viewed and subsequently visualized carefully parametrized visual
73 objects. First, we found that, as in non-human primates (57–60), human VTC neurons showed
74 robust visual responses (61–65), and are well modelled by linearly combining features in deep
75 layers of a deep network trained to perform object classification (66). We confirmed two
76 consequences of such a code: each neuron had a linear null space orthogonal to its encoding axis
77 (given by the coefficients of linear combination); furthermore, each neuron responded maximally
78 to synthetic stimuli generated using its encoding axis. Second, we asked subjects to imagine (i.e.,
79 visualize) a diverse subset of objects that they had previously seen, while recording responses of
80 visually-characterized VTC neurons. We found that a subset (~40%) of axis-tuned neurons
81 reactivated during imagery, and the imagined responses of individual neurons to specific objects
82 were proportional to the projection value of those objects onto the neurons' preferred axes.

83 Together these findings provide evidence for the implementation of a generative model in human
84 VTC by neurons that represent both real and imagined stimuli.
85

86 **Results**

87
88 **Neurons show diverse category tuning**
89
90 We examined how human VTC neurons encode visual objects by recording responses of these
91 neurons to a large set of objects with varying features using a rapid screening task. Patients
92 sequentially viewed a series of 500 images (4 repetitions each, 2000 total trials), drawn from face,
93 text, plant, animal, and object categories (Figure 1A, top; Figure S1A for detailed schematic; Figure
94 S2A for stimulus examples). At random intervals a catch question pertaining to the immediately
95 preceding image would appear on the screen. Images stayed on screen for 250 ms and the inter
96 trial interval was randomized between 100-150 ms. Despite the rapid presentation rate, patients
97 answered 77% of the catch questions correctly, indicating that the stimuli were carefully attended
98 to (Figure 1B).
99

100 We recorded 714 VTC neurons in 57 sessions across 16 patients. Response onset latency of
101 individual neurons was computed on a trial-by-trial basis using a Poisson burst metric (67) (see
102 methods). Out of the 714 neurons, 456 showed a significantly increased response to the onset of
103 a visual stimulus (Figure 1C; 1 x 5 sliding window ANOVA, bin size 50 ms, step size 5 ms, 5
104 consecutive significant bins with $p < 0.01$; see methods). The locations of all electrodes from
105 which we recorded visually responsive neurons can be seen in Figure 1D (red and blue dots);
106 supplementary tables S1 and S2 contain neuron count and patient demographic information. The
107 mean response latency was 162 ms (Figure 1E). Human VTC neurons showed diverse tuning
108 patterns, including neurons that were selectively silent to a given category (Figure 1F), neurons
109 that were maximally responsive to a given category (Figure 1G, top), neurons characterized by an
110 initial suppression of activity (Figure 1G, middle), and neurons that distinguished categories via a
111 latency code (Figure 1G, bottom).
112

113 **Neurons are tuned to specific axes in object space**

114
115 We next examined whether visually responsive neurons encoded specific object features. We
116 leveraged deep networks trained on object classification to build a low-dimensional object space
117 that captures the shape and appearance of arbitrary objects (68) without relying on subjective
118 visual descriptions. We built our object space by passing the 500 images shown to patients
119 through AlexNet (66) and performing principal components analysis (PCA) on the unit activations
120 in layer fc6 (Figure 2A). We determined that 50 dimensions explained 80.68% of the variance in
121 fc6 responses (Figure S2B), and used these 50 dimensions in all remaining analyses. This allowed
122 us to describe every visual object shown to patients as a point in a 50-dimensional feature space.
123 We next investigated how neural activity mapped onto this feature space.
124

125 For each cell, we computed a ‘preferred axis’ given by the coefficients c in the equation $\vec{r} = c \cdot$
126 $\vec{f} + \vec{c}_0$, where \vec{r} is the response of the neuron to a given image, \vec{f} is the 50D object feature vector
127 of that image, and \vec{c}_0 is a constant offset (see methods). The axis tuning of an example neuron is
128 shown in Figure 2B. This neuron showed a monotonically increasing response to stimuli with
129 higher projection values onto its preferred axis (see methods) while showing no change in tuning
130 along an axis orthogonal to its preferred axis. We confirmed that the correlation between
131 projection value and firing rate was significant when compared to a shuffled distribution along
132 the preferred axis ($p < 0.01$, Figure 2B, top distribution) and not significant along the principal
133 orthogonal axis ($p = 0.496$, Figure 2B, bottom distribution). This “axis code” emphasizes the
134 geometric picture that neurons are projecting incoming stimuli, formatted as vectors in the
135 complex feature space, onto their specific preferred axes in the space (57, 58, 69).

136

137 Axis tuning also clarified the response pattern of some neurons with complex tuning. The example
138 neuron in Figure 2C responded most strongly to the plant category, but also responded strongly
139 to specific stimuli in other categories (see raster, Figure 2C). Indeed, no semantic label could
140 obviously delineate between the stimuli that elicited the strongest and weakest response from
141 the neuron (Figure 2D). However, this neuron showed significant axis tuning (Figure 2E), and
142 sorting the stimuli according to their projection value onto the preferred axis shows that the top
143 stimuli were those with the highest projection values and vice versa (Figure 2F). Across the
144 population, we found that a majority (367/456, ~80%) of visually responsive neurons were
145 significantly axis tuned (Figure 2G, H), showing a significant positive correlation between
146 projection value and response along the preferred axis with no such correlation along the
147 orthogonal axis (Figure 2I; $r_{pref} = 0.54$, $r_{ortho} = -1.3833e-10$, $p = 1.43e-121$, Wilcoxon ranksum
148 test).

149

150 We subsequently compared the variance explained by the axis model to that explained by two
151 alternative models: a category label model and an exemplar model. The category label model was
152 chosen for comparison due to the robust category responses seen in VTC (Figure 1F, G). The
153 exemplar model is a well-known alternative to the axis model (70–72) which posits that neurons
154 have maximal responses to specific exemplars (i.e., specific points in object space) and decaying
155 responses to objects with increasing distance from the exemplar (see methods). Only 18/456
156 neurons had good fits (positive explained variance) across all models (34/456 exemplar model,
157 343/456 category label, 339/456 axis model) and even in these neurons the axis model explained
158 significantly more variance than either alternative (39.62% axis, 17.04% category, $p = 9.42e-04$,
159 Wilcoxon ranksum test; 39.62% axis, 10.37% exemplar, $p = 4.78e-05$, Wilcoxon ranksum test;
160 Figure 2J, K).

161

162 We built our low-dimensional object space by leveraging AlexNet, a deep network trained to
163 perform object classification. However, there now exist a plethora of deep convolutional neural
164 network models that are capable of performing object recognition at or beyond human capability
165 (73). A comparison of several such models (VGG-16/19/Face and CORNet models; see
166 supplementary methods) revealed that axis tuning in feature spaces built from the fully
167 connected layers of these models explained a large amount of the explainable neural variance
168 (~45%) in the 367 axis-tuned VTC neurons (Figure S3, S4). Moreover, the proportion of explainable

169 variance explained was roughly the same for all models except VGG-Face and an ‘eigen-model’
170 where PCA is conducted directly on the pixels; implying that the axis code is independent of the
171 specific convolutional network used to build the feature space.
172

173 Reconstructing objects using human VTC responses

174
175 We next investigated the richness of the VTC representation by attempting to reconstruct objects
176 using VTC responses. A consequence of the linear relationship between VTC neuron responses
177 and object features is the ease of learning a linear decoder that predicts object feature values
178 from the population activity (74–76). The responses of the neurons can be approximated as a
179 linear combination of object features, with the slopes of the ramps corresponding to the weights
180 (69–71). Then for a population of neurons, $\vec{R} = C\vec{F} + \vec{C}_0$, where \vec{R} is the response vector of the
181 different neurons, C is the weight matrix, \vec{F} is the vector of object feature values, and \vec{C}_0 is the
182 offset vector. Inverting this equation to solve for \vec{F} provides a linear decoder that predicts object
183 feature values from the population activity (74–76) (see methods). We used this approach
184 coupled with leave-one-out cross-validation to learn the linear transform that maps responses to
185 features (Figure 3A).
186

187 Figure 3B shows the correlation between actual feature values and model predictions for the first
188 two dimensions of object space; both correlations were significant positive values ($r_{\text{dim } 1} = 0.46$,
189 $p = 2.6\text{e-}28$, paired t-test, $r_{\text{dim } 2} = 0.69$, $p = 0$, paired t-test); Figure S2C shows correlations for all
190 50 dimensions. To reconstruct objects, we searched a large auxiliary object database for the
191 object with the feature vector closest to that decoded from the neural activity. A normalized
192 distance in feature space between the best possible and actual reconstruction was computed to
193 quantify the decoding accuracy across all stimuli (Figure 3C, see methods). Side-by-side
194 comparisons of the original images with the ‘reconstructions’ chosen from the auxiliary database
195 showed striking visual correspondence (Figure 3D). We also confirmed that VTC responses could
196 support well-above-chance decoding of the correct object from distractors (Figure 3E).
197

198 Predicting responses to synthetic stimuli

199
200 To further validate the axis model, we attempted to predict the responses of VTC neurons to
201 images not shown to our patients, in a “closed loop” experiment. We used an initial screening
202 session to compute axes for all neurons recorded. We then used a generative adversarial network
203 (GAN) trained to invert the responses of AlexNet (77, 78) to systematically generate images
204 corresponding to evenly spaced points along both the preferred and orthogonal axes. We then
205 returned to the patient room for a second session and rescreened with the synthetic images
206 added to the original stimulus set. We performed this experiment in 8 sessions across 4 patients
207 which yielded 16 axis tuned neurons.
208

209 We predicted that images with projection values greater than the maximum projection value of
210 the original stimulus images should serve as ‘super-stimuli’ driving the neuron to a higher
211 response than any of the original stimulus images (78). Figure 3F shows an example neuron tested

212 with the synthetic stimuli. This neuron's most preferred stimulus was a t-shirt, while the least
213 preferred stimulus was a ladder (Figure 3F); these two stimuli were also the ones that projected
214 most and least strongly onto the cell's preferred axis. Synthetic stimuli were sampled evenly along
215 both preferred (yellow dots) and orthogonal axes (green dots) such that the extremes had larger
216 projection values than any of the original stimulus images (Figure 3F). Responses to these
217 synthetic stimuli demonstrate the expected increase along the preferred axis and no increase
218 along the orthogonal axis (Spearman's rank correlation between projection value and firing rate
219 for preferred axis 0.8303, $p = 4e-3$ when compared to shuffled distribution, orthogonal axis = -
220 0.1037, $p = 0.604$; Figure 3F). Moreover, the firing rates to the subset of synthetic stimuli with
221 larger projection values along the preferred axis than the original stimulus images were
222 substantially higher than those to the original stimulus images. Interestingly, the most and least
223 effective synthetic stimuli showed striking resemblance to the most and least effective original
224 stimuli. Figure 3G shows responses of the same cell to synthetic stimuli evenly sampled in a 2D
225 grid spanned by the preferred and orthogonal axes (purple square, Figure 3F); responses showed
226 the expected changes in tuning only along directions parallel to the preferred axis. Figure 3H
227 shows a second example neuron which was face selective ($r_{pref} = 0.89$, $p = 5.53e-04$; $r_{ortho} =$
228 0.52, $p = 0.12$).
229

230 Across all 16 neurons the correlation between projection value and firing rate for the synthetic
231 images along both was significantly higher along the preferred axis (Figure 3I, $r_{pref} = 0.58$, r_{ortho}
232 = 0.08, $p = 1.34e-05$, Wilcoxon ranksum test). In this experiment the axes were computed using
233 the 500 original stimulus images only (see methods). The axes were also used to predict the firing
234 rate responses to the synthetic stimuli. The distribution of correlation values between the
235 responses predicted using the axes computed in the first session and the recorded responses in
236 the second session are shown in Figure 3J (mean = 0.58).
237

238 Neuronal activity during imagery

239 So far, we have established that human VTC uses an axis code to represent objects and confirmed
240 this code through stimulus reconstruction and generation of super-stimuli. Armed with this
241 understanding of the sensory code, we could now tackle the single neuron mechanisms of visual
242 imagery. In 12 sessions across 6 patients, we tested patients on a cued imagery task following
243 visual screening (Figure 1A, bottom; Figure S1B for detailed schematic). In the imagery task,
244 patients viewed and subsequently visualized from memory 6-8 objects out of the original 500
245 used for screening. Each trial required alternating, cued visualization of two object stimuli after
246 an initial encoding period during which patients passively viewed the two stimuli. During the cued
247 imagery period patients would close their eyes and imagine the 2 objects in the trial in alternating
248 5 s periods until each stimulus had been imagined 4 times, being informed to switch images at
249 the end of a given 5 s period verbally by the experimenter. Each image appeared in 2 separate
250 trials for a total of 8 imagery trials per image (see methods). We recorded 231 VTC neurons in the
251 imagery task, of which 131 were visually responsive and 107 were axis tuned.
252

253

254 We found robust activation of neurons in human VTC during imagery, with 66/231 neurons
255 (~30%) showing activation to at least one object (Figure 4A, 1 x n sliding window ANOVA or sliding
256 window ttest, n = number of stimuli, bin size 1.5 s, step size 300 ms, 6 consecutive bins with p <
257 0.05 – see methods). As in vision, neurons showed a diverse range of response profiles during
258 imagery, some activating sparsely during imagery of a single specific stimulus (Figure 4B), others
259 activating to imagery of multiple stimuli in a graded manner (Figure 4C), and a small number
260 (15/231) activating only during imagery and not during viewing (Figure 4D).

261

262 VTC neurons recapitulate the visual code during imagery

263

264 A central goal of the current study was to clarify whether imagery reactivates visual neurons in a
265 way that respects their perceptual code, or whether an alternative code (possibly implemented
266 by a distinct population of neurons) is recruited during imagery. The former would constitute
267 strong evidence for the existence of a generative model in the human brain. To establish whether
268 VTC neurons reactivate in a manner that respects the axis code, the rapid screening session was
269 conducted using the 500 object images described earlier and axes were computed for the axis-
270 tuned neurons. Then, 6-8 stimuli that were spread along the preferred and orthogonal axes were
271 chosen for use in the cued imagery task. Figures 5A&B show two example neurons: both showed
272 significant axis tuning and responded most strongly to the image with the largest projection value
273 onto the preferred axes during both encoding and recall (piano and mirror respectively).

274

275 Examining the population of axis tuned neurons that reactivated during imagery (43/107, ~40%)
276 revealed a significant correlation between projection value onto the neurons' preferred axes
277 (computed using screening responses) and responses during imagery ($r_{pref} = 0.20$, Figure 5E left,
278 $p = 0.001$ as compared to a shuffled distribution, Figure 5F top) with no such correlation along
279 the orthogonal axes ($r_{ortho} = -0.08$, Figure 5E right, $p = 0.929$, Figure 5F bottom). Lastly, the
280 distribution of Spearman's rank correlation coefficients between viewing and imagery of the same
281 images in reactivated axis tuned neurons showed a high median value significantly larger than 0
282 (0.18 ± 0.05 , $p = 2.4e-03$, one sample ttest; Figure 5G). Taken together these findings indicate that
283 neurons in human VTC support visual imagery by reinstating similar network activity to viewing.

284

285 Discussion

286

287 In this paper, we explore the long-standing hypothesis that mental imagery is supported by
288 reactivation of sensory areas in the visual domain. We focused our investigations on VTC, a region
289 long known to harbor representations of complex visual objects (57, 58, 64). We first mapped out
290 the feedforward code for visual objects, and then measured responses of the same population
291 during imagery. We find that, as in non-human primates (57, 58, 69), human VTC neurons (367
292 out of 456 visually responsive neurons recorded) represented visual objects via linear projection
293 of incoming object vectors onto a specific 'preferred' axis in a high dimensional feature space
294 built using the unit activations of a deep network. Confirming the axis model, we could
295 reconstruct viewed objects with high accuracy using a linear decoder (Figure 3A-E), and generate

296 synthetic super-stimuli for cells using the axis mapped to real stimuli (Figure 3F-J). These results
297 indicate that the VTC neurons reported in this study serve to recognize objects by simply
298 measuring their features not identifying them semantically, and that axis tuning is a powerful
299 quantitative way to conceptualize the response of a substantial portion of human VTC neurons.
300 Thus it would appear that in macaque, human, and deep neural network (DNN) architectures, an
301 essential stage of object processing relies on a meaning-agnostic distributed shape
302 representation (57–60, 68, 79).

303

304 A subset of neurons (66/231 total, 43/107 axis tuned) reactivated during imagery in a manner
305 that respected the axis code: imagined responses were significantly correlated to projection value
306 onto the preferred axis but not the orthogonal axis (Figure 5F), and viewed and imagined
307 responses were positively correlated (Figure 5G). No previous single neuron study has examined
308 visual imagery while recording from neurons in human VTC whose sensory code was
309 characterized in detail. A single neuron study of imagery in the medial temporal lobe and another
310 of spoken free recall (a related behavior) in human VTC demonstrated reactivation of a few
311 neurons in both areas (62, 80), but in both of these studies the sensory code was unknown. Our
312 findings demonstrate for the first time that neurons in VTC support visual imagery by reactivating
313 in a structured manner that respects the visual code.

314

315 The source of the top-down signal driving VTC reactivation during imagery remains an open
316 question. Candidates for this source include the hippocampus and prefrontal cortex, given their
317 involvement in various forms of memory (81–87), their dense connections to VTC (88, 89), and
318 the known ability of human hippocampal neurons to be selectively reactivated by free recall (19,
319 90). Another question is the relationship between the VTC signals during imagery and those
320 previously reported in primary visual areas (V1/V2) (39, 91, 92). Given hierarchically organized
321 feedback connections (88), we hypothesize that the VTC signals may be driving the imagery-
322 related signals in earlier visual areas. Future work is required to investigate the response
323 characteristics of reactivated neurons across the brain, including those both upstream and
324 downstream of VTC during imagery.

325

326 This work is the first to reveal a detailed understanding of the neural codes underlying visual
327 object perception and imagery in human VTC. In particular, the results provide evidence for the
328 existence of a generative model in the human brain—a mechanism capable of synthesizing
329 detailed sensory contents from an abstract, semantic representation (93–95), effectively inverting
330 the classic feedforward pathway. Generative models derive incredible computational power by
331 transforming challenges in perception and cognition, such as inference under ambiguity (52, 93,
332 94) and object-based attention (94), into closed-loop feedback systems (96). The existence of a
333 generative model in the human brain may even explain creative artistic processes that have so far
334 remained out of reach of neuroscientific understanding.

335

336

337

338

339 **Methods**

340

341 **Participants**

342

343 The study participants were 16 adult patients who were implanted with depth electrodes for
344 seizure monitoring as part of an evaluation for treatment of drug-resistant epilepsy (see Table S2
345 for demographic data). All patients provided informed consent and volunteered to participate in
346 this study. Research protocols were approved by the institutional review board of Cedars-Sinai
347 Medical Center (Study 572). The tasks were conducted while patients stayed in the epilepsy
348 monitoring unit following implantation of depth electrodes. The location of the implanted
349 electrodes was solely determined by clinical needs. The neural results were analyzed across all 16
350 patients. Each of the 16 patients included in this study had at least one depth electrode targeting
351 the ventral temporal cortex.

352

353

354 **Psychophysical tasks**

355

356 Patients participated in three different tasks: an initial screening, cued imagery, and a final re-
357 screening. The initial screening session was conducted to identify axis tuned neurons. Then 6-8
358 stimuli were chosen for the imagery task that had some spread along both the preferred and
359 orthogonal axes. After axis tuned neurons were identified and the stimuli chosen we then
360 conducted the cued imagery task, followed by another screening session immediately after. The
361 second screening was used to match the neurons from the first and second sessions.

362

363 **Screening**

364

365 Patients viewed a set of 500 object stimuli (grayscale, white background, size: 224x224) with
366 varying features (taken from www.freepngs.com) 4 times each for a total of 2000 trials in a
367 shuffled order. Images were displayed on laptop computer with a 15.5" screen placed 1 meter
368 away and subtended 6-7 visual degrees. Each image stayed on screen for 250 ms and the inter-
369 trial interval consisting of a blank screen was jittered between 100-150 ms. The task was
370 punctuated with yes/no 'catch' questions pertaining to the image that came right before the
371 question requiring the patients to pay close attention in order to answer them correctly (Figure
372 S1A). Catch questions occurred between 2 to 80 trials after the previous one. Patients responded
373 to the questions using an RB-844 response pad (https://cedrus.com/rb_series/). During the
374 synthetic image screens, the synthetic images would be added to the original stimulus set and
375 the task parameters remained unchanged.

376

377 **Cued Imagery**

378

379 Patients viewed a set of 6-8 object stimuli chosen from the 500 used for screening (taken from
380 www.freepngs.com). Each trial focused on 2 images and had an encoding period, a visual search

381 distraction period, and a cued imagery period. During encoding patients would see the 2 images
382 4 times each in a shuffled order. Each image stayed on screen for 1.5 s and the inter-image interval
383 was 1.5-2 s. After this encoding period a visual search puzzle was presented (puzzles created by
384 artist Gergely Dudas <https://thedudolf.blogspot.com/>) and stayed on screen for 30 s. After
385 reporting via button press whether or not they were able to find the object in the puzzle, patients
386 began the cued imagery period. During cued imagery patients would close their eyes and imagine
387 the stimuli in the trial in an alternating fashion for 4 repetitions of 5 s each (40 s continuous
388 imagery period). Patients were cued to switch the image they were imagining every 5 s by verbal
389 cue (Figure S1B). After the imagery period patients would begin the next trial via button press
390 when ready. Every image was present in 2 trials, leading to 8 repetitions of both encoding and
391 imagery for each image.

392

393

394 Electrophysiology

395

396 The data in this paper was recorded from left and/or right VTC in addition to the other clinically
397 relevant targets (unique to each patient) using Behnke-Fried micro-macro electrodes (Ad-Tech
398 Medical Instrument Corporation) (95). All analyses in this paper are based on the signals recorded
399 from the 8 microwires protruding from the end of the electrode. Recordings were performed with
400 an FDA-approved electrophysiology system, and sampled at 32khz (ATLAS, Neuralynx Inc.) (97).

401

402 Spike sorting and quality metrics

403

404 Signals were bandpass filtered offline in the range of 300-3000hz with a zero phase lag filter
405 before spike detection. Spike detection and sorting was carried out via the semiautomated
406 template matching algorithm Osort (98). The properties of clusters identified as putative neurons
407 and subsequently used for analysis were documented using a suite of spike quality metrics (Figure
408 S5).

409

410 Electrode localization

411

412 Electrode localization was based on postoperative imaging using either MRI or computed
413 tomography (CT) scans. We co-registered postoperative and preoperative MRIs using Freesurfer's
414 mri_robust_register (99). To summarize recording locations across participants we aligned each
415 participant's preoperative MRI to the CITI168 template brain in MNI152 coordinates (100) using
416 a concatenation of an affine transformation and symmetric image normalization (SyN)
417 diffeomorphic transform (101). The MNI coordinates of the microwires from a given electrode
418 shank were marked as one location. MNI coordinates of microwires with putative neurons
419 detected from all participants were plotted on a template brain for visualization (Figure 1D).

420

421

422

423 Data Analyses

424

425 Visual responsiveness classification

426

427 To assess whether a neuron was visually responsive we used a 1 x 5 sliding window ANOVA with
428 the factor visual category (face, plant, animal, text, object). We counted spikes in an 80-400 ms
429 period relative to stimulus onset, using a bin size of 50 ms and a step size of 5 ms. Beginning at
430 each time point, the average response in each 50 ms trial snippet was computed and the vector
431 of responses, labeled by their stimulus identity, were fed into the ANOVA. The time point was
432 then incremented by 5 ms and the ANOVA was re-computed. A neuron was considered visually
433 responsive if the ANOVA was significant ($p < 0.01$) for 6 consecutive time points. These
434 parameters were chosen to ensure that the probability of selecting a neuron by chance was less
435 than 0.05 (compared to bootstrap distribution with 1000 repetitions, Figure S6A).

436

437 Response latency computation

438

439 We computed such a single trial onset latency by using a Poisson spike-train analysis for all visually
440 responsive neurons. This method detects points of time in which the observed inter-spike
441 intervals (ISI) deviate significantly from that assumed by a constant-rate Poisson process. This is
442 done by maximizing a Poisson surprise index (83). The mean firing rate of the neuron during the
443 inter-trial interval was used to set the baseline rate for the Poisson process. Spikes from a window
444 of 80-300 ms after stimulus onset were included. For a given burst of spikes, if the probability that
445 said burst was produced by a constant-rate Poisson process - where the rate parameters are
446 specified by the baseline firing rate - was less than 0.001, we took the timepoint of the first spike
447 as the onset latency. The response latency of the neuron was taken to be the average latency
448 across all trials.

449

450 Building an object space using a deep network

451

452 We built a high dimensional object space by feeding our 500 stimulus images into the pre-trained
453 MATLAB implementation of AlexNet (90) (Deep Learning Toolbox, command: 'net = alexnet'). The
454 responses of the 4096 nodes in fc6 were extracted to form a 500 x 4096 matrix (using the
455 'activations' function). PCA was then performed on this matrix yielding 499 PCs, each of length
456 4096. To reduce the dimensionality of this space we retained only the first 50 PCs which captured
457 80.68% of the response variance across fc6 units (Figure S2B). The first two dimensions accounted
458 for 20.17% of the response variance across fc6 units.

459

460 Axis computation

461

462 Preferred axis

463

464 The preferred axis of each neuron was computed using the spike triggered average (STA). The
465 neural response vector was computed by binning spikes elicited by each stimulus in a 250ms
466 window starting from the response latency of the neuron — necessarily restricting analysis to
467 visually responsive neurons.

468

469 Once the neural response vector was computed the STA is defined as:

470

471
$$\vec{P}_{sta} = (\vec{r} - \bar{r})^T F \quad (1)$$

472

473 where \vec{r} is the $n \times 1$ neural response vector to n objects, \bar{r} is the mean firing rate, and F is an $n \times$
474 d matrix of features with each row corresponding to the features for a given object that is
475 computed via PCA on deep network activations (see above). The projection value of the stimulus
476 objects onto the preferred axis is given by:

477

478
$$Proj_{sta} = \left(\frac{\vec{P}_{sta}}{\|\vec{P}_{sta}\|} \right) F^T \quad (2)$$

479

480 Principal orthogonal axis

481

482 The orthogonal axis seen in all plots is the principal orthogonal axis. This is defined as the axis
483 orthogonal to the preferred axis along which there is the most variation. For each neuron, the
484 preferred axis was computed, the component along the preferred axis (\vec{P}) was subsequently
485 subtracted from all object feature vectors in F leaving a matrix of orthogonal feature vectors.
486 Succinctly, for a given feature vector \vec{f}_d in feature space we computed:

487

488
$$\vec{f}_{d-1} = \vec{f}_d - \vec{f}_d^T \vec{P} \frac{\vec{P}}{\|\vec{P}\|^2} \quad (3)$$

489

490 Then principal component analysis was performed on this set of n vectors \vec{f}_{d-1} , and the first
491 principal component is chosen as the principal orthogonal axis.

492

493 Quantifying significance of axis tuning

494

495 For each neuron after the preferred axis was computed we examined the correlation between
496 the firing rate response to the stimuli and their projection value along the preferred axis. This
497 correlation value was recomputed after shuffling the features (1000 repetitions) and the original
498 value was compared to this bootstrap distribution. If the original value was greater than 99% of
499 the shuffled values the neuron was considered axis tuned.

500

501 Explained variance computation

502

503 Axis model

504
505 The axis model assumes a linear relationship between the projection value of an incoming object
506 onto the neurons preferred axis and its response. Therefore, to quantify the explained variance
507 for each neuron, we fit a linear regression model between the PCs of the features and the
508 responses of the neuron. A leave-one-out cross validation approach was used i.e. the responses
509 to 499 objects were used to fit the model, and the responses of the neuron to the left-out object
510 was predicted using the same linear transform. In this manner we could produce a predicted
511 response for all images. Note that the computation of variance explained by the category label
512 was done in this manner as well, replacing the PC features with the vector of category labels.
513

514 **Exemplar model**

515
516 The exemplar model assumes that each neuron has a maximal response to a specific exemplar in
517 object space and that the response of the neuron to an incoming object decays as a function of
518 the distance from the object to this exemplar. We used a previous implementation of the
519 exemplar model (92) in which the response of the neuron which has an exemplar \vec{e} to an incoming
520 object \vec{x} is:

521
522
$$f(d) = C_0 + C_1 d + C_2 d^2 + C_3 d^3 \quad (4)$$

523

524
525 where d is the Euclidean distance between the exemplar object and the incoming object:
526

527
528
$$d = \sqrt{\sum_{i=1}^N (e_i - x_i)^2} \quad (5)$$

529 and N is the dimensionality of the object space, which in our analyses was 50. In such an
530 implementation the coefficients of the polynomial C and the features of the exemplar \vec{e} are
531 considered free parameters. They were adjusted iteratively to minimize the error of fit using the
532 MATLAB function *lsqcurvefit*.

533
534 To set an upper bound for the explained variance different trials of responses to the stimuli were
535 randomly split into two halves. The Pearson correlation (r) between the average responses from
536 two half-splits across images was calculated and corrected using the Spearman-Brown correction:
537

538
539
$$r' = \frac{2r}{(1+r)} \quad (6)$$

540 The square of r' was considered the upper bound or explainable variance. The reported results
541 are the ratio of explained to explainable variance.
542

543 **Decoding analysis**

544
545 We find that neurons in human VTC are performing linear projection onto specific preferred axis
546 in object space. As such, their responses can be well modeled as a population by the equation:
547

548
$$\vec{R} = [C \vec{C}_o] \begin{bmatrix} \vec{F} \\ 1 \end{bmatrix} \quad (7)$$

549
550 where \vec{R} is the $n \times 1$ population response vector to a given image (n = number of neurons), C is
551 the $n \times d$ weight matrix for different neurons (d = number of dimensions i.e. 50), \vec{F} is the $d \times 1$
552 vector of object feature values, and \vec{C}_o is the $n \times 1$ offset vector. Thus, the decoding analysis was
553 performed by inverting (7), yielding:

554
555
$$\begin{bmatrix} \vec{F} \\ 1 \end{bmatrix} = [C \vec{C}_o]^+ \vec{R} = K \vec{R} \quad (8)$$

556
557 Where C^+ indicates the Moore-Penrose pseudoinverse, and K is the $d + 1 \times n$ matrix that
558 transforms measured firing rates \vec{R} into predicted features \vec{F} . We used the responses of all but
559 one of the objects ($500 - 1 = 499$) to determine K using the MATLAB function *regress*. These were
560 then substituted into (8) to predict the feature vector of the last object. For the i^{th} image

561
562
$$\vec{F}_i = \sum_{d=1}^{N_{\text{dim}}} (k_{o_d} + \sum_{n=1}^{N_{\text{cells}}} r_{i,n} k_{n,d}) \quad (9)$$

563
564 Where $r_{i,n}$ is the n^{th} element of \vec{R} , i.e. the response of neuron n to the image i , $k_{n,d}$ is the (n, d)
565 element of the weight (regression coefficient) matrix, and k_{o_d} is the d^{th} element of the offset
566 vector. Decoding accuracy was quantified by randomly selecting a subset of object images that
567 included the actual feature vector of the decoded object from the total set of 500 and compared
568 their feature vectors to the predicted feature vector of the decoded object by Euclidean distance.
569 If the actual feature vector closest to the predicted feature vector is of the object being decoded
570 ('target') the decoding is considered correct. This procedure is repeated 1000 times for each of
571 the 500 images with a varying number of distractors to get an aggregate measure of decoding
572 accuracy (Figure 3E).

573
574 **Object 'reconstruction'**

575
576 To generate images that reflect the features encoded in the neural responses we gathered images
577 from an auxiliary database and passed 15,901 background free images through AlexNet. The
578 images were then projected into the space built by the 500 stimulus objects. None of these ~16k
579 images had been shown to the patients. For each stimulus image the feature vector decoded from
580 the neural activity was compared to the feature vectors of the large stimulus set. The object in
581 the large image set with the smallest Euclidean distance to the decoded feature vector was
582 considered the 'reconstruction' of that stimulus image (94).

583
584 To account for the fact that the large object set did not contain any images shown to the patients,
585 which sets a limit on how good the reconstruction can be, we computed a ‘normalized distance’
586 to quantify the reconstruction accuracy for each object. We defined the normalized
587 reconstruction distance for an image as

588

$$589 \text{Normalized distance} = \frac{|V_{\text{recon}} - V_{\text{original}}|}{|V_{\text{best possible recon}} - V_{\text{original}}|} \quad (10)$$

590
591 where V_{recon} is the feature vector reconstructed from neuronal responses, V_{original} is the feature
592 vector of the image presented to the patients, and $V_{\text{best possible recon}}$ is the feature vector of the
593 best possible reconstruction (image in the large set with the closest distance to V_{original}). A
594 normalized distance of 1 means the decoded image is the best reconstruction possible. The
595 median normalized distance value for our data was 2.256. For a vast majority (482/500, ~96%) of
596 the images the normalized distance was small (< 5, Figure 3C), with the distance of the worst
597 performing image being 10.471, implying that the neural responses captured many of the fine
598 feature details of the original objects.

599
600 **Generation of synthetic stimuli**

601
602 The axis model provides a very clear relationship between images and responses for individual
603 neurons. In essence, images with increasing projection values onto a neuron’s preferred axis will
604 show increasing firing rates. This implies that if one computes a neuron’s preferred axis and then
605 evenly samples points along it and generates images from those points, those images will elicit
606 systematically increasing responses from the neuron. This also implies that if one generates an
607 image from a point further along the axis than any of the stimulus images used to compute the
608 neurons axis, that image will act as a super stimulus and drive the neuron to a higher firing rate
609 than any of the stimulus images.

610
611 To test these predictions we ran the screening task in one session, computed the axes for the
612 neurons recorded, sampled points along the preferred and orthogonal axes and fed those vectors
613 back into a pre-trained GAN (90) to generate the synthetic stimuli. We then went back to the
614 patient room and re-ran the screening task with the synthetic images added in. The neurons from
615 the first and second sessions were matched (see below) and the responses of the neurons to the
616 synthetic stimuli were recorded.

617
618 **Computation of predicted responses to synthetic images**

619
620 Predicted responses to the synthetic images were computed by fitting a linear regression model
621 between the PCs of the features and the responses of the neuron during the first session to the
622 500 stimulus images. That linear transform was then used to predict the responses to the
623 synthetic images. These predicted responses were then compared to the responses recorded in
624 matching neuron during the second session.

625

626 Matching neurons across experiment sessions

627

628 Given the nature of recording neurons in a clinical setting wherein you can only record a few
629 neurons at a time it is common practice to run an initial screening task in order to determine the
630 stimuli that drive the neurons being recorded before using those stimuli in other tasks. In such
631 cases it is generally assumed that the neurons recorded a few hours apart are the same but it is
632 important to provide evidence.

633

634 One method of verification is to re-run the same screen in a subsequent session and compare the
635 selectivity of the neurons in both sessions. However, if there are multiple neurons having roughly
636 similar selectivity (which was sometimes the case in our data) matching individual neurons is
637 difficult. In order to meet this challenge, we examined multiple features of the neurons recorded
638 in the first and second sessions. Our algorithm would compute the selectivity vector (rank ordered
639 list of stimulus number for the neuron), the waveform, the burst index which is a measure of how
640 many bursts per unit time the neuron discharges (98), the computed response latency of the
641 neuron, and whether or not the neuron was axis-tuned (binary variable). The selectivity vectors
642 were compared using cosine distance (MATLAB function *pdist*), and the waveforms by Euclidean
643 distance.

644

645 Each axis-tuned neuron in the initial session was compared to all ipsilateral axis-tuned neurons in
646 the subsequent session (same procedure for non-axis tuned neurons). The session two neurons
647 were then rank ordered in every category with first place in a given category giving a session two
648 neuron a score of x where $x = \text{number of session two neurons being compared to the one session}$
649 $\text{one neuron, second place receiving a score of } x - 1 \text{ all the way until the last place neuron in a}$
650 $\text{given category receives a score of 1. The scores of all session two neurons were then summed}$
651 $\text{and the algorithm would assign the session two neuron with the max score to be the session one}$
652 $\text{neuron's 'match'. This procedure is then repeated in the reverse direction, i.e. each session two}$
653 $\text{neuron is compared to all session one neurons. Pairs that were bijective were automatically}$
654 $\text{returned as 'matches' and all others were marked out for manual curation. Manual curation was}$
655 $\text{carried out by examining the used metrics in addition to shape of the peri-stimulus time}$
656 $\text{histogram of the category response of all potential neuron match pairs.}$

657

658 Reactivation metric

659

660 To assess whether a neuron was active during imagery we used a combination of a $1 \times N$ ($N =$
661 $\text{number of images})$ sliding window ANOVA and a sliding window ttest during the cued imagery
662 period. We counted spikes in a 0-5 s window relative to imagery onset (i.e. the entire cued
663 imagery period), using a bin size of 1.5 s and a step size of 300 ms. Beginning at each time point,
664 the trial average was computed using spikes in a 1.5 s window and the vector of trials was fed
665 into the ANOVA or ttest. For the ANOVA the trials were labeled by their stimulus identity. The
666 time point was incremented by 300 ms and the ANOVA was re-computed. A neuron was
667 considered active during imagery if either the ANOVA or the ttest was significant ($p < 0.05$) for 6

668 consecutive bins (or 5 consecutive steps). These parameters were chosen to ensure that the
669 probability of selecting a neuron by chance ($p_{\text{ANOVA}} + p_{\text{ttest}}$) was less than 0.05 (compared to
670 bootstrap distribution with 1000 repetitions, Figure S6B&C).

671

672 Correlation of viewed and imagined responses

673

674 For neurons that were active during imagery we computed the correlation between viewed and
675 imagined responses. To compute the viewed response to each stimulus, we collected spikes in a
676 1 s window of the encoding period starting from the response latency of the neuron (computed
677 during screening) and averaged across repetitions of a given stimulus. To compute the imagined
678 response for all neurons active during imagery, we collected spikes in a 2 s window starting from
679 the first significant time bin (1 x N sliding window ANOVA or sliding window ttest, N = number of
680 stimuli, $p < 0.05$) and averaged across repetitions of a given stimulus. The Spearman rank
681 correlation (r) was then computed between these 2 vectors.

682

683

684

685

686

687

688

689 **Acknowledgements:** We thank the staff of the epilepsy monitoring unit (neuromonitoring
690 staff, nursing staff, and physicians) and of the Biomedical Imaging Research Institute at Cedars-
691 Sinai Medical Center for patient care and support. We thank Emily Choe for help implementing
692 the GAN used in Figure 3. We thank members of the Tsao and Rutishauser labs namely Janis
693 Hesse, Hristos Courellis, and Francesco Lanfranchi for helpful comments throughout all stages of
694 this project. We thank the patients for all their patience and perseverance.

695

696 **Funding:** This work was funded by the BRAIN initiative through the NIH Office of the Director
697 (U01NS117839), the Howard Hughes Medical Institute, the Simons Foundation Collaboration on
698 the Global Brain, and the Chen Center for Systems Neuroscience at Caltech.

699

700 **Author contributions:** V.W., U.R., and D.Y.T designed the study. V.W. collected and analyzed
701 the data. V.W., U.R., and D.Y.T. wrote the paper with input from C.M.R, J.M.C., L.M.B., and A. N. M.
702 C.M.R, J.M.C., and L.M.B. provided patient care and facilitated experiments. A. N. M. performed
703 surgery.

704

705 **Competing interests:** Authors declare no competing interests.

706

707 **Data and materials availability:** Data and code will be made publicly available upon
708 acceptance.

709 **References**

- 710
- 711 1. A. Schlegel, P. Alexander, S. V. Fogelson, X. Li, Z. Lu, P. J. Kohler, E. Riley, P. U. Tse, M. Meng,
712 The artist emerges: visual art learning alters neural structure and function. *Neuroimage* **105**,
713 440–451 (2015).
- 714 2. J. F. Christensen, A. Gomila, Introduction: Art and the brain: From pleasure to well-being.
715 *Prog. Brain Res.* **237**, xxvii–xlvi (2018).
- 716 3. M. Ellamil, C. Dobson, M. Beeman, K. Christoff, Evaluative and generative modes of thought
717 during the creative process. *Neuroimage* **59**, 1783–1794 (2012).
- 718 4. S. S. H. Wong, S. W. H. Lim, Mental imagery boosts music compositional creativity. *PLoS One*
719 **12**, e0174009 (2017).
- 720 5. L. Taruffi, M. B. Küssner, A review of music-evoked visual mental imagery: Conceptual issues,
721 relation to emotion, and functional outcome. *Psychomusicology: Music, Mind, and Brain* **29**,
722 62–74 (2019).
- 723 6. S. Liu, H. M. Chow, Y. Xu, M. G. Erkkinen, K. E. Swett, M. W. Eagle, D. A. Rizik-Baer, A. R. Braun,
724 Neural correlates of lyrical improvisation: an fMRI study of freestyle rap. *Sci. Rep.* **2**, 834
725 (2012).
- 726 7. P. E. Keller, Mental imagery in music performance: underlying mechanisms and potential
727 benefits. *Ann. N. Y. Acad. Sci.* **1252**, 206–213 (2012).
- 728 8. K. Erhard, F. Kessler, N. Neumann, H.-J. Ortheil, M. Lotze, Professional training in creative
729 writing is associated with enhanced fronto-striatal activity in a literary text continuation task.
730 *Neuroimage* **100**, 15–23 (2014).
- 731 9. C. Shah, K. Erhard, H.-J. Ortheil, E. Kaza, C. Kessler, M. Lotze, Neural correlates of creative
732 writing: an fMRI study. *Hum. Brain Mapp.* **34**, 1088–1101 (2013).
- 733 10. S. Liu, M. G. Erkkinen, M. L. Healey, Y. Xu, K. E. Swett, H. M. Chow, A. R. Braun, Brain activity
734 and connectivity during poetry composition: Toward a multidimensional model of the
735 creative process. *Hum. Brain Mapp.* **36**, 3351–3372 (2015).
- 736 11. D. L. Schacter, D. R. Addis, R. L. Buckner, Episodic simulation of future events: concepts, data,
737 and applications. *Ann. N. Y. Acad. Sci.* **1124**, 39–60 (2008).
- 738 12. S. T. Moulton, S. M. Kosslyn, Imagining predictions: mental imagery as mental emulation.
739 *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **364**, 1273–1280 (2009).

- 740 13. C. Guariglia, L. Pizzamiglio, "The Role of Imagery in Navigation: Neuropsychological
741 Evidence" in *Spatial Processing in Navigation, Imagery and Perception*, F. Mast, L. Jäncke,
742 Eds. (Springer US, Boston, MA, 2007), pp. 17–28.
- 743 14. C. M. Bird, J. A. Bisby, N. Burgess, The hippocampus and spatial constraints on mental
744 imagery. *Front. Hum. Neurosci.* **6**, 142 (2012).
- 745 15. A. Bocchi, M. Carrieri, S. Lancia, V. Quaresima, L. Piccardi, The Key of the Maze: The role of
746 mental imagery and cognitive flexibility in navigational planning. *Neurosci. Lett.* **651**, 146–
747 150 (2017).
- 748 16. E. Tulving, Others, Episodic and semantic memory. *Organization of memory* **1**, 1 (1972).
- 749 17. D. R. Addis, A. T. Wong, D. L. Schacter, Remembering the past and imagining the future:
750 common and distinct neural substrates during event construction and elaboration.
751 *Neuropsychologia* **45**, 1363–1377 (2007).
- 752 18. D. L. Schacter, D. R. Addis, R. L. Buckner, Remembering the past to imagine the future: the
753 prospective brain. *Nat. Rev. Neurosci.* **8**, 657–661 (2007).
- 754 19. H. Gelbard-Sagiv, R. Mukamel, M. Harel, R. Malach, I. Fried, Internally generated reactivation
755 of single neurons in human hippocampus during free recall. *Science* **322**, 96–101 (2008).
- 756 20. A. P. Vaz, J. H. Wittig Jr, S. K. Inati, K. A. Zaghloul, Replay of cortical spiking sequences during
757 human memory retrieval. *Science* **367**, 1131–1134 (2020).
- 758 21. A. Mary, J. Dayan, G. Leone, C. Postel, F. Fraisse, C. Malle, T. Vallée, C. Klein-Peschanski, F.
759 Viader, V. de la Sayette, D. Peschanski, F. Eustache, P. Gagnepain, Resilience after trauma:
760 The role of memory suppression. *Science* **367** (2020).
- 761 22. T. D. Griffiths, Musical hallucinosis in acquired deafness. Phenomenology and brain
762 substrate. *Brain* **123** (Pt 10), 2065–2076 (2000).
- 763 23. S. Kumar, W. Sedley, G. R. Barnes, S. Teki, K. J. Friston, T. D. Griffiths, Neural bases of musical
764 hallucinations. *J. Neurol. Neurosurg. Psychiatry* **85**, e3–e3 (2014).
- 765 24. R. J. Zatorre, A. R. Halpern, D. W. Perry, E. Meyer, A. C. Evans, Hearing in the mind's ear: A
766 PET investigation of musical imagery and perception. *J. Cogn. Neurosci.* **8**, 29–46 (1996).
- 767 25. A. R. Halpern, R. J. Zatorre, When that tune runs through your head: a PET investigation of
768 auditory imagery for familiar melodies. *Cereb. Cortex* **9**, 697–704 (1999).
- 769 26. R. J. Stevenson, T. I. Case, Olfactory imagery: a review. *Psychon. Bull. Rev.* **12**, 244–264 (2005).
- 770 27. J. A. Gottfried, A. P. R. Smith, M. D. Rugg, R. J. Dolan, Remembrance of odors past. *Neuron*
771 **42**, 687–695 (2004).

- 772 28. S.-S. Yoo, D. K. Freeman, J. J. McCarthy III, F. A. Jolesz, Neural substrates of tactile imagery: a
773 functional MRI study. *Neuroreport* **14**, 581–585 (2003).
- 774 29. X. Tian, J. M. Zarate, D. Poeppel, Mental imagery of speech implicates two mechanisms of
775 perceptual reactivation. *Cortex* **77**, 1–12 (2016).
- 776 30. W. Richter, R. Somorjai, R. Summers, M. Jarmasz, R. S. Menon, J. S. Gati, A. P. Georgopoulos,
777 C. Tegeler, K. Ugurbil, S. G. Kim, Motor area activity during mental rotation studied by time-
778 resolved single-trial fMRI. *J. Cogn. Neurosci.* **12**, 310–320 (2000).
- 779 31. R. N. Shepard, J. Metzler, Mental rotation of three-dimensional objects. *Science* **171**, 701–
780 703 (1971).
- 781 32. L. M. Parsons, P. T. Fox, J. H. Downs, T. Glass, T. B. Hirsch, C. C. Martin, P. A. Jerabek, J. L.
782 Lancaster, Use of implicit motor imagery for visual shape discrimination as revealed by PET.
783 *Nature* **375**, 54–58 (1995).
- 784 33. J. Decety, The neurophysiological basis of motor imagery. *Behav. Brain Res.* **77**, 45–52 (1996).
- 785 34. K. J. Miller, G. Schalk, E. E. Fetz, M. den Nijs, J. G. Ojemann, R. P. N. Rao, Cortical activity
786 during motor execution, motor imagery, and imagery-based online feedback. *Proc. Natl.
787 Acad. Sci. U. S. A.* **107**, 4430–4435 (2010).
- 788 35. S. M. Kosslyn, A. Pascual-Leone, O. Felician, S. Camposano, J. P. Keenan, W. L. Thompson, G.
789 Ganis, K. E. Sukel, N. M. Alpert, The role of area 17 in visual imagery: convergent evidence
790 from PET and rTMS. *Science* **284**, 167–170 (1999).
- 791 36. A. Spagna, D. Hajhajate, J. Liu, P. Bartolomeo, Visual mental imagery engages the left fusiform
792 gyrus, but not the early visual cortex: A meta-analysis of neuroimaging evidence. *Neurosci.
793 Biobehav. Rev.* **122**, 201–217 (2021).
- 794 37. J. Pearson, T. Naselaris, E. A. Holmes, S. M. Kosslyn, Mental Imagery: Functional Mechanisms
795 and Clinical Applications. *Trends Cogn. Sci.* **19**, 590–602 (2015).
- 796 38. M. B. Bone, M. St-Laurent, C. Dang, D. A. McQuiggan, J. D. Ryan, B. R. Buchsbaum, Eye
797 Movement Reinstatement and Neural Reactivation During Mental Imagery. *Cereb. Cortex* **29**,
798 1075–1089 (2019).
- 799 39. S. D. Slotnick, W. L. Thompson, S. M. Kosslyn, Visual mental imagery induces retinotopically
800 organized activation of early visual areas. *Cereb. Cortex* **15**, 1570–1583 (2005).
- 801 40. S.-H. Lee, D. J. Kravitz, C. I. Baker, Disentangling visual imagery and perception of real-world
802 objects. *Neuroimage* **59**, 4064–4073 (2012).
- 803 41. A. Ishai, J. V. Haxby, L. G. Ungerleider, Visual imagery of famous faces: effects of memory and
804 attention revealed by fMRI. *Neuroimage* **17**, 1729–1741 (2002).

- 805 42. N. Dijkstra, S. E. Bosch, M. A. J. van Gerven, Shared Neural Mechanisms of Visual Perception
806 and Imagery. *Trends Cogn. Sci.* **23**, 423–434 (2019).
- 807 43. A. Ishai, D. Sagi, Common mechanisms of visual imagery and perception. *Science* **268**, 1772–
808 1774 (1995).
- 809 44. C. Hofstetter, A. Achaibou, P. Vuilleumier, Reactivation of visual cortex during memory
810 retrieval: content specificity and emotional modulation. *Neuroimage* **60**, 1734–1745 (2012).
- 811 45. K. M. O’Craven, N. Kanwisher, Mental imagery of faces and places activates corresponding
812 stimulus-specific brain regions. *J. Cogn. Neurosci.* **12**, 1013–1023 (2000).
- 813 46. J. Pearson, The human imagination: the cognitive neuroscience of visual mental imagery.
814 *Nat. Rev. Neurosci.* **20**, 624–634 (2019).
- 815 47. L. G. Ungerleider, J. V. Haxby, “What” and “where” in the human brain. *Curr. Opin. Neurobiol.*
816 **4**, 157–165 (1994).
- 817 48. M. J. Farah, D. N. Levine, R. Calvanio, A case study of mental imagery deficit. *Brain Cogn.* **8**,
818 147–164 (1988).
- 819 49. M. J. Riddoch, Loss of visual imagery: A generation deficit. *Cogn. Neuropsychol.* **7**, 249–273
820 (1990).
- 821 50. S. Thorudottir, H. M. Sigurdardottir, G. E. Rice, S. J. Kerry, R. J. Robotham, A. P. Leff, R. Starrfelt,
822 The architect who lost the ability to imagine: The cerebral basis of visual imagery. *Brain Sci.*
823 **10**, 59 (2020).
- 824 51. R. P. Rao, D. H. Ballard, Predictive coding in the visual cortex: a functional interpretation of
825 some extra-classical receptive-field effects. *Nat. Neurosci.* **2**, 79–87 (1999).
- 826 52. Y. Huang, J. Gornet, S. Dai, Z. Yu, T. Nguyen, D. Y. Tsao, A. Anandkumar, Neural networks with
827 recurrent generative feedback, *arXiv [cs.LG]* (2020). <http://arxiv.org/abs/2007.09200>.
- 828 53. P. Cavanagh, G. P. Caplovitz, T. K. Lytchenko, M. Maechler, P. U. Tse, D. Sheinberg, Object-
829 Based Attention, *PsyArXiv* (2022). <https://doi.org/10.31234/osf.io/2bsn7>.
- 830 54. C. G. Gross, C. E. Rocha-Miranda, D. B. Bender, Visual properties of neurons in inferotemporal
831 cortex of the Macaque. *J. Neurophysiol.* **35**, 96–111 (1972).
- 832 55. R. Desimone, T. D. Albright, C. G. Gross, C. Bruce, Stimulus-selective properties of inferior
833 temporal neurons in the macaque. *J. Neurosci.* **4**, 2051–2062 (1984).
- 834 56. I. Fried, C. L. Wilson, N. T. Maidment, J. Engel Jr, E. Behnke, T. A. Fields, K. A. MacDonald, J.
835 W. Morrow, L. Ackerson, Cerebral microdialysis combined with single-neuron and

- 836 electroencephalographic recording in neurosurgical patients. Technical note. *J. Neurosurg.*
837 **91**, 697–705 (1999).
- 838 57. L. Chang, D. Y. Tsao, The Code for Facial Identity in the Primate Brain. *Cell* **169**, 1013–1028.e14
839 (2017).
- 840 58. P. Bao, L. She, M. McGill, D. Y. Tsao, A map of object space in primate inferotemporal cortex.
841 *Nature* **583**, 103–108 (2020).
- 842 59. Y. Shi, D. Bi, J. K. Hesse, F. F. Lanfranchi, S. Chen, D. Y. Tsao, Rapid, concerted switching of the
843 neural code in inferotemporal cortex. *bioRxiv.org*, 2023.12. 06.570341 (2023).
- 844 60. L. She, M. K. Benna, Y. Shi, S. Fusi, D. Y. Tsao, Temporal multiplexing of perception and
845 memory codes in IT cortex. *Nature* **629**, 861–868 (2024).
- 846 61. V. Axelrod, C. Rozier, T. S. Malkinson, K. Lehongre, C. Adam, V. Lambrecq, V. Navarro, L.
847 Naccache, Face-selective neurons in the vicinity of the human fusiform face area. *Neurology*
848 **92**, 197–198 (2019).
- 849 62. S. Khuvivis, E. M. Yeagle, Y. Norman, S. Grossman, R. Malach, A. D. Mehta, Face-Selective Units
850 in Human Ventral Temporal Cortex Reactivate during Free Recall. *J. Neurosci.* **41**, 3386–3399
851 (2021).
- 852 63. V. Axelrod, C. Rozier, T. S. Malkinson, K. Lehongre, C. Adam, V. Lambrecq, V. Navarro, L.
853 Naccache, Face-selective multi-unit activity in the proximity of the FFA modulated by facial
854 expression stimuli. *Neuropsychologia* **170**, 108228 (2022).
- 855 64. R. Quijan Quiroga, M. Boscaglia, J. Jonas, H. G. Rey, X. Yan, L. Maillard, S. Colnat-Coulbois, L.
856 Koessler, B. Rossion, Single neuron responses underlying face recognition in the human
857 midfusiform face-selective cortex. *Nat. Commun.* **14**, 5661 (2023).
- 858 65. H. S. Courellis, J. Minxha, A. R. Cardenas, D. L. Kimmel, C. M. Reed, T. A. Valiante, C. D.
859 Salzman, A. N. Mamelak, S. Fusi, U. Rutishauser, Abstract representations emerge in human
860 hippocampal neurons during inference. *Nature* **632**, 841–849 (2024).
- 861 66. A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional
862 neural networks. *Adv. Neural Inf. Process. Syst.* **25** (2012).
- 863 67. D. P. Hanes, K. G. Thompson, J. D. Schall, Relationship of presaccadic activity in frontal eye
864 field and supplementary eye field to saccade initiation in macaque: Poisson spike train
865 analysis. *Exp. Brain Res.* **103**, 85–96 (1995).
- 866 68. D. L. K. Yamins, H. Hong, C. F. Cadieu, E. A. Solomon, D. Seibert, J. J. DiCarlo, Performance-
867 optimized hierarchical models predict neural responses in higher visual cortex. *Proc. Natl.*
868 *Acad. Sci. U. S. A.* **111**, 8619–8624 (2014).

- 869 69. L. She, M. K. Benna, Y. Shi, S. Fusi, D. Y. Tsao, The neural code for face memory. *BioRxiv* (2021).
- 870 70. T. Valentine, A unified account of the effects of distinctiveness, inversion, and race in face
871 recognition. *Q. J. Exp. Psychol. A* **43**, 161–204 (1991).
- 872 71. M. Riesenhuber, T. Poggio, Hierarchical models of object recognition in cortex. *Nat. Neurosci.*
873 **2**, 1019–1025 (1999).
- 874 72. D. Y. Tsao, W. A. Freiwald, What's so special about the average face? *Trends Cogn. Sci.* **10**,
875 391–393 (2006).
- 876 73. Y. Taigman, M. Yang, M. Ranzato, L. Wolf, “DeepFace: Closing the gap to human-level
877 performance in face verification” in *2014 IEEE Conference on Computer Vision and Pattern
878 Recognition* (IEEE, 2014), pp. 1701–1708.
- 879 74. K. N. Kay, T. Naselaris, R. J. Prenger, J. L. Gallant, Identifying natural images from human brain
880 activity. *Nature* **452**, 352–355 (2008).
- 881 75. A. S. Cowen, M. M. Chun, B. A. Kuhl, Neural portraits of perception: reconstructing face
882 images from evoked brain activity. *Neuroimage* **94**, 12–22 (2014).
- 883 76. A. Nestor, D. C. Plaut, M. Behrmann, Feature-based face representations and image
884 reconstruction from behavioral and neural data. *Proc. Natl. Acad. Sci. U. S. A.* **113**, 416–421
885 (2016).
- 886 77. A. Dosovitskiy, T. Brox, Generating images with perceptual similarity metrics based on deep
887 networks. *Adv. Neural Inf. Process. Syst.* **29** (2016).
- 888 78. C. R. Ponce, W. Xiao, P. F. Schade, T. S. Hartmann, G. Kreiman, M. S. Livingstone, Evolving
889 Images for Visual Neurons Using a Deep Generative Network Reveals Coding Principles and
890 Neuronal Preferences. *Cell* **177**, 999-1009.e10 (2019).
- 891 79. J. Kubilius, M. Schrimpf, A. Nayebi, D. Bear, D. L. K. Yamins, J. J. DiCarlo, CORnet: Modeling
892 the Neural Mechanisms of Core Object Recognition, *bioRxiv* (2018)p. 408385.
- 893 80. G. Kreiman, C. Koch, I. Fried, Imagery neurons in the human brain. *Nature* **408**, 357–361
894 (2000).
- 895 81. J. Daume, J. Kaminski, A. G. P. Schjetnan, Y. Salimpour, U. Khan, C. Reed, W. Anderson, T. A.
896 Valiante, A. N. Mamelak, U. Rutishauser, Control of working memory maintenance by theta-
897 gamma phase amplitude coupling of human hippocampal neurons. *bioRxiv.org*, doi:
898 10.1101/2023.04.05.535772 (2023).
- 899 82. J. Kamiński, S. Sullivan, J. M. Chung, I. B. Ross, A. N. Mamelak, U. Rutishauser, Persistently
900 active neurons in human medial frontal and medial temporal lobe support working memory.
901 *Nat. Neurosci.* **20**, 590–601 (2017).

- 902 83. J. Kamiński, A. Brzezicka, A. N. Mamelak, U. Rutishauser, Combined phase-rate coding by
903 persistently active neurons as a mechanism for maintaining multiple items in working
904 memory in humans. *Neuron* **106**, 256–264.e3 (2020).
- 905 84. E. Svoboda, M. C. McKinnon, B. Levine, The functional neuroanatomy of autobiographical
906 memory: a meta-analysis. *Neuropsychologia* **44**, 2189–2208 (2006).
- 907 85. U. Rutishauser, I. B. Ross, A. N. Mamelak, E. M. Schuman, Human memory strength is
908 predicted by theta-frequency phase-locking of single neurons. *Nature* **464**, 903–907 (2010).
- 909 86. W. B. Scoville, B. Milner, Loss of recent memory after bilateral hippocampal lesions. *J. Neurol.*
910 *Neurosurg. Psychiatry* **20**, 11–21 (1957).
- 911 87. H. Tomita, M. Ohbayashi, K. Nakahara, I. Hasegawa, Y. Miyashita, Top-down signal from
912 prefrontal cortex in executive control of memory retrieval. *Nature* **401**, 699–703 (1999).
- 913 88. D. J. Felleman, D. C. Van Essen, Distributed hierarchical processing in the primate cerebral
914 cortex. *Cereb. Cortex* **1**, 1–47 (1991).
- 915 89. P. Grimaldi, K. S. Saleem, D. Tsao, Anatomical connections of the functionally defined “face
916 patches” in the macaque monkey. *Neuron* **90**, 1325–1342 (2016).
- 917 90. M. Cerf, N. Thiruvengadam, F. Mormann, A. Kraskov, R. Q. Quiroga, C. Koch, I. Fried, On-line,
918 voluntary control of human temporal lobe neurons. *Nature* **467**, 1104–1108 (2010).
- 919 91. B. Thirion, E. Duchesnay, E. Hubbard, J. Dubois, J.-B. Poline, D. Lebihan, S. Dehaene, Inverse
920 retinotopy: inferring the visual content of images from brain activation patterns.
921 *Neuroimage* **33**, 1104–1116 (2006).
- 922 92. T. Naselaris, C. A. Olman, D. E. Stansbury, K. Ugurbil, J. L. Gallant, A voxel-wise encoding
923 model for early visual areas decodes mental images of remembered scenes. *Neuroimage*
924 **105**, 215–228 (2015).
- 925 93. T. S. Lee, D. Mumford, Hierarchical Bayesian inference in the visual cortex. *J. Opt. Soc. Am. A*
926 *Opt. Image Sci. Vis.* **20**, 1434–1448 (2003).
- 927 94. A. Yuille, D. Kersten, Vision as Bayesian inference: analysis by synthesis? *Trends Cogn. Sci.* **10**,
928 301–308 (2006).
- 929 95. B. Peters, J. J. DiCarlo, T. Gureckis, R. Haefner, L. Isik, J. Tenenbaum, T. Konkle, T. Naselaris, K.
930 Stachenfeld, Z. Tavares, D. Tsao, I. Yildirim, N. Kriegeskorte, How does the primate brain
931 combine generative and discriminative computations in vision?, *arXiv [q-bio.NC]* (2024).
932 <http://arxiv.org/abs/2401.06005>.
- 933 96. Y. Ma, D. Tsao, H.-Y. Shum, On the principles of Parsimony and Self-consistency for the
934 emergence of intelligence, *arXiv [cs.AI]* (2022). <http://arxiv.org/abs/2207.04630>.

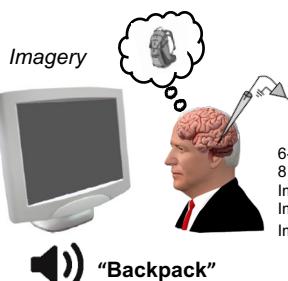
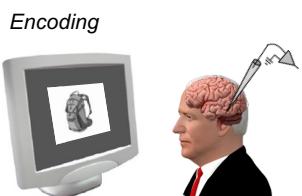
- 935 97. J. Minxha, A. N. Mamelak, U. Rutishauser, Surgical and electrophysiological techniques for
936 single-neuron recordings in human epilepsy patients. *Extracellular recording* (2018).
- 937 98. U. Rutishauser, E. M. Schuman, A. N. Mamelak, Online detection and sorting of
938 extracellularly recorded action potentials in human medial temporal lobe recordings, *in vivo*.
939 *J. Neurosci. Methods* **154**, 204–224 (2006).
- 940 99. M. Reuter, H. D. Rosas, B. Fischl, Highly accurate inverse consistent registration: a robust
941 approach. *Neuroimage* **53**, 1181–1196 (2010).
- 942 100. W. M. Pauli, A. N. Nili, J. M. Tyszka, A high-resolution probabilistic *in vivo* atlas of human
943 subcortical brain nuclei. *Sci Data* **5**, 180063 (2018).
- 944 101. B. Avants, J. T. Duda, J. Kim, H. Zhang, J. Pluta, J. C. Gee, J. Whyte, Multivariate analysis of
945 structural and diffusion imaging in traumatic brain injury. *Acad. Radiol.* **15**, 1360–1375
946 (2008).
- 947 102. A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, I. Sutskever, Zero-shot
948 text-to-image generation. *ICML* **abs/2102.12092**, 8821–8831 (2021).
- 949
- 950
- 951
- 952
- 953
- 954
- 955
- 956
- 957
- 958
- 959
- 960
- 961
- 962
- 963
- 964

A



500 Stimuli
4 repetitions
Image ON: 250ms
Image OFF: 100-150ms

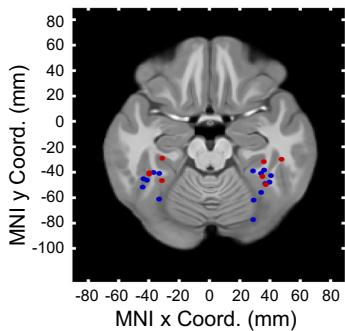
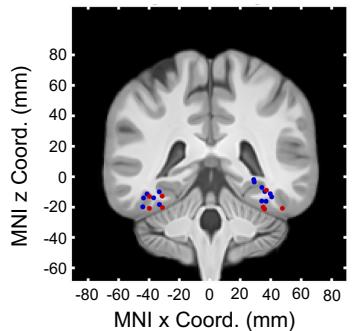
Cued Imagery Task



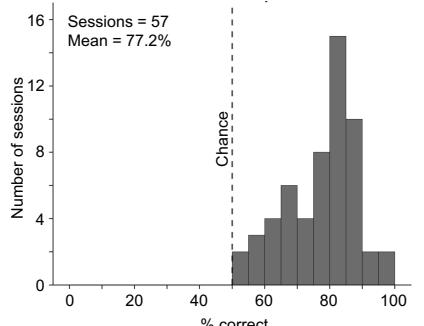
6-8 Stimuli
8 repetitions
Image ON: 1.5s
Image OFF: 1.5-2s
Imagery ON: 5s

D

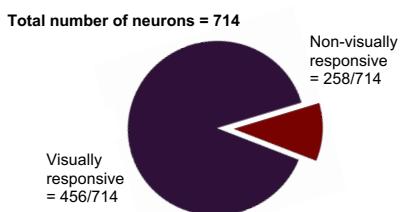
- Screening
- Screening & cued imagery



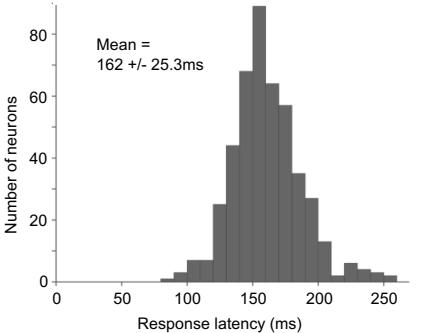
B



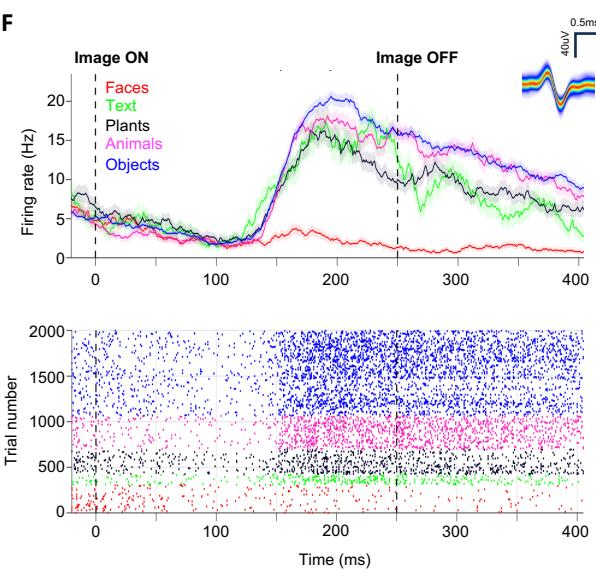
C



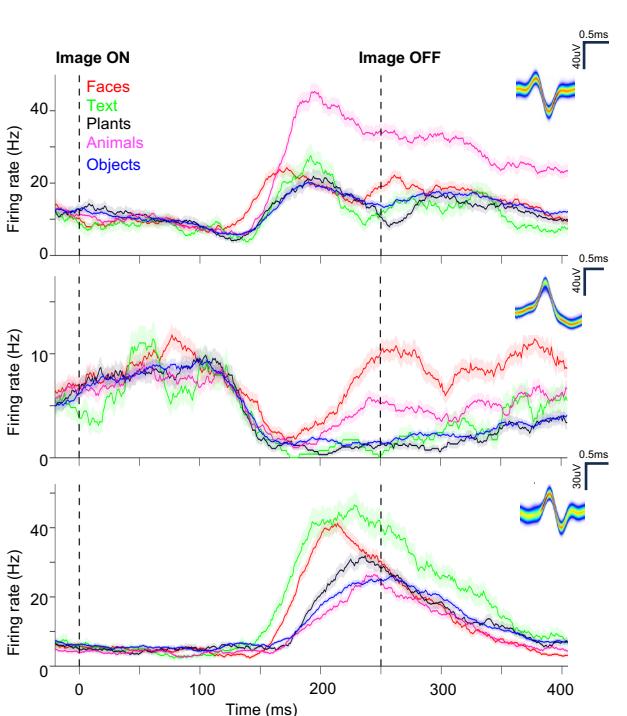
E



F



G



965 **Figure 1. Selective visual responses in human VTC.**

966 **(A)** Task schematics. Patients performed two tasks: a screening task (top) and a cued imagery task
967 (bottom). In the screening task, grayscale images with white backgrounds were displayed on a
968 gray screen for 250 ms with the inter-trial interval jittered between 100-150 ms. Images
969 subtended 6-7 visual degrees. At random intervals (min interval: 1 trial, max interval: 80 trials) a
970 yes-no catch question would appear pertaining to the image that came just before it. **(B)** Accuracy
971 of catch question responses for all sessions. On average patients answered catch questions
972 correctly in $77 \pm 11\%$ of the trials (\pm s.d.) despite the rapid stimulus presentation, implying the
973 stimuli were closely attended to. **(C)** 456 of the 714 recorded neurons were visually responsive.
974 **(D)** Recording locations of the 27 microwire bundles (left and right) that contained at least one
975 well isolated neuron in human VTC across all 16 patients in Coronal (left) and Axial (right) views.
976 Montreal Neurological Institute coordinates can be read off the image or seen in Table S1. Each
977 dot represents the location of one microwire bundle (8 channels). The locations marked in red
978 were sessions used in a subsequent cued imagery task. **(E)** Distribution of response latencies for
979 all 456 visually responsive VTC neurons. The mean response latency was 162 ± 25 ms (\pm s.d.). **(F)**
980 An example neuron recorded during the screening task. This neuron was responsive to all
981 categories except faces. The estimated response latency is 178 ms. Stimulus onset is at $t = 0$ and
982 offset is at $t = 250$ ms. The inset shows the mean waveform of the neuron. **(G)** Further example
983 neurons illustrating the diversity of response profiles. (Top) A strong category selective neuron,
984 showing a lower response to any of the non-preferred categories. (Middle) A response profile
985 characterized by an initial suppression of activity. (Bottom) A neuron that distinguishes its
986 preferred from non-preferred category using a latency code along with a rate code.

987

988

989

990

991

992

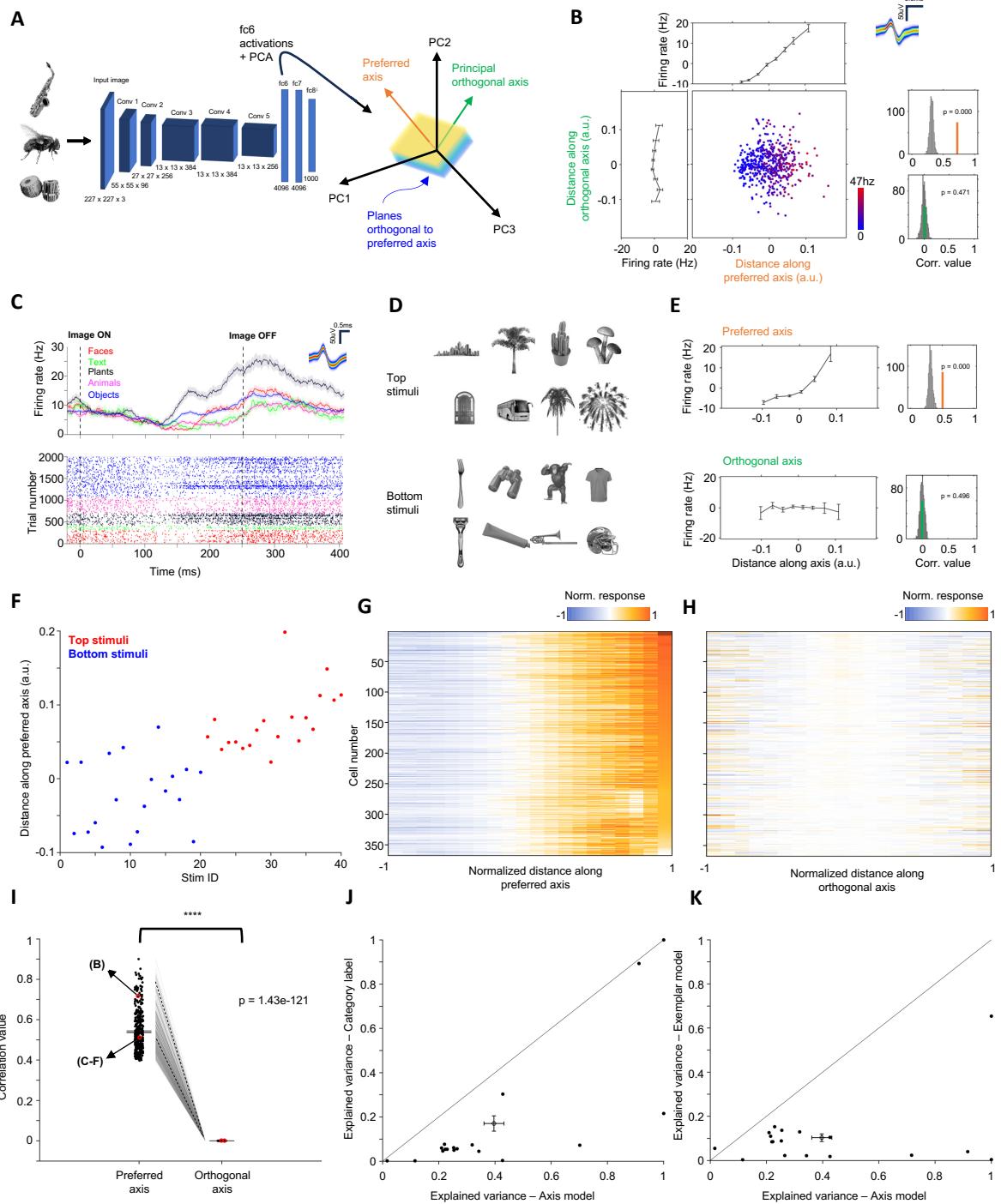
993

994

995

996

997



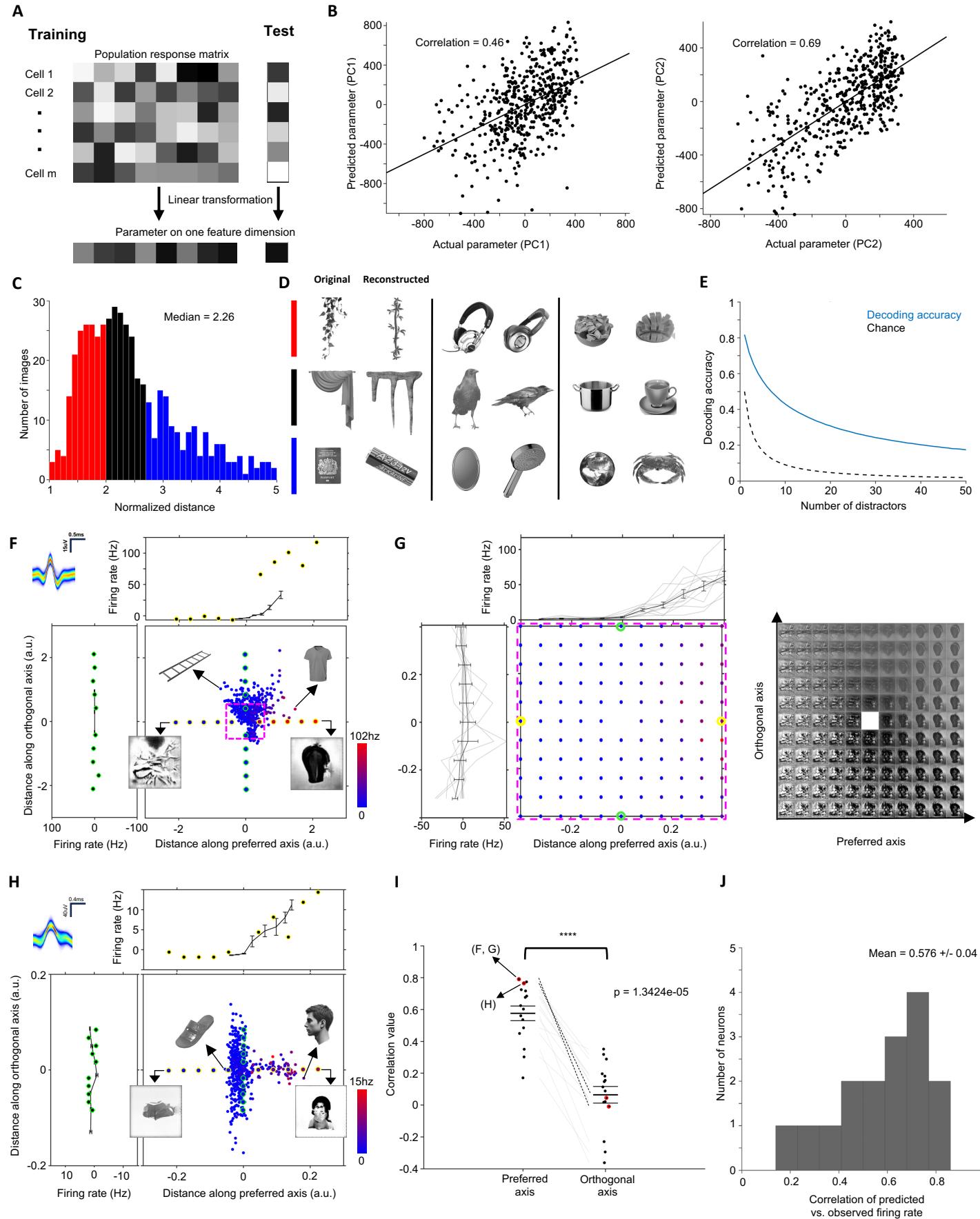
998 **Figure 2. Axis tuning in human VTC neurons.**

999 **(A)** Schematic of the stimulus parametrization and axis computation procedure. Individual stimuli
1000 were parametrized as points in a 50-dimensional feature space that was built by performing PCA
1001 on the unit activations of AlexNet's fc6 layer and keeping the top 50 PCs. Neurons' preferred axes
1002 were computed in this space by projecting the features onto the mean subtracted responses
1003 leading to a large null space for each neuron. **(B)** Example axis-tuned neuron. Scatter shows the
1004 responses to the 500 stimulus images projected onto the neuron's preferred axis and principal
1005 (longest) orthogonal axis in the face feature space. Response magnitude is color-coded. (Top)
1006 Mean response as a function of distance along the preferred axis. (Left) Mean response as a
1007 function of distance along the orthogonal axis. (Top distribution) The correlation value between
1008 projection value and firing rate for the preferred axis (orange) showing a significant correlation (p
1009 = 0.001, bootstrap distribution 1000 repetitions). (Bottom distribution) Correlation between the
1010 projection value and firing along the orthogonal axis (green) is not significantly different from the
1011 null distribution for the orthogonal axis (p = 0.496, bootstrap distribution 1000 repetitions). **(C-F)**
1012 The axis model explains complex neural tuning that does not follow pre-defined categorical
1013 boundaries. **(C)** An example neuron is shown. Peri-stimulus time histogram (PSTH) and raster of
1014 the example neuron. A close look at the raster reveals robust responses to a few out of
1015 "preferred" category objects. **(D)** The top (most preferred) and bottom (least preferred) stimuli
1016 for the neuron shown in (C), no semantic category cleanly delineates between them. **(E)** Axis
1017 tuning for the neuron shown in (C-D). See (B) for notation. **(F)** Projection values of the stimuli
1018 shown in (D) for the neuron shown in (C-D) reveals a systematic relationship between the
1019 projection value and stimulus preference. **(G-J)** Population summary of the tuning to the
1020 preferred (G) and orthogonal (H) axes for all visually responsive neurons that were significantly
1021 axis tuned (367/456). Each row is a neuron and shows the increase in normalized response as a
1022 function of distance along the preferred axis. **(H)** Corresponding plot for the principal orthogonal
1023 axis (orthogonal axis capturing the most variation) showing no systematic change in response as
1024 a function of distance. **(I)** Pearson correlation between the projection value of the stimulus images
1025 onto the preferred and orthogonal axes and the firing rate response of the neuron ($n=367$). The
1026 neurons shown in (B) and (C-F) are marked in red. **(J)** Comparison of fit quality for axis model and
1027 category label for 18 neurons. These 18 neurons were chosen as they had significant variance
1028 explained in their responses across all 3 models. The axis model provides significantly better fits
1029 to actual responses than the category label. (39.62% axis, 17.04% category, p = 9.42e-04,
1030 Wilcoxon ranksum test). The error bars indicate standard error of the mean. **(K)** Comparison of fit
1031 quality for axis and exemplar models for the same 18 neurons. The axis model provides
1032 significantly better fits to actual responses than the exemplar model (39.62% axis, 10.37%
1033 exemplar, p = 4.78e-05, Wilcoxon ranksum test). The error bars indicate standard error of the
1034 mean.

1035

1036

1037



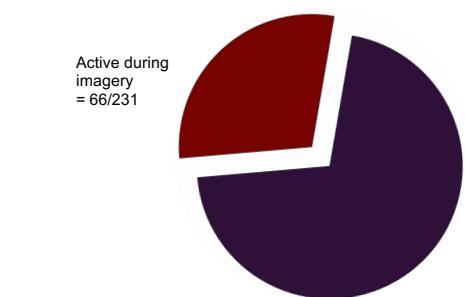
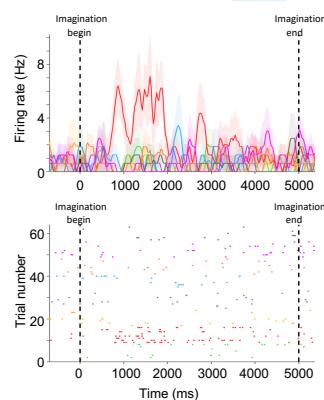
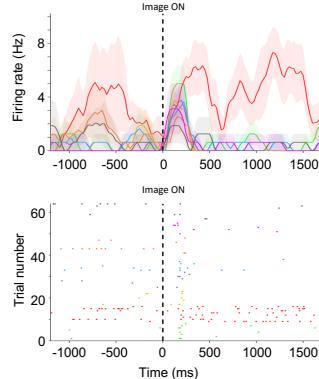
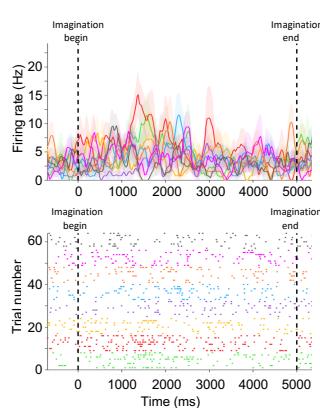
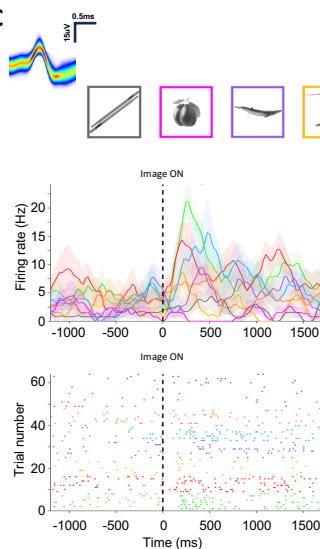
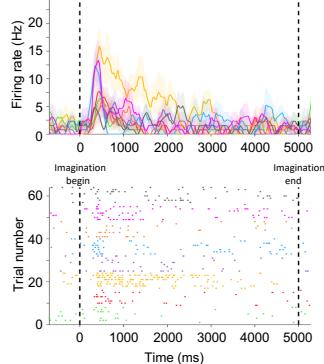
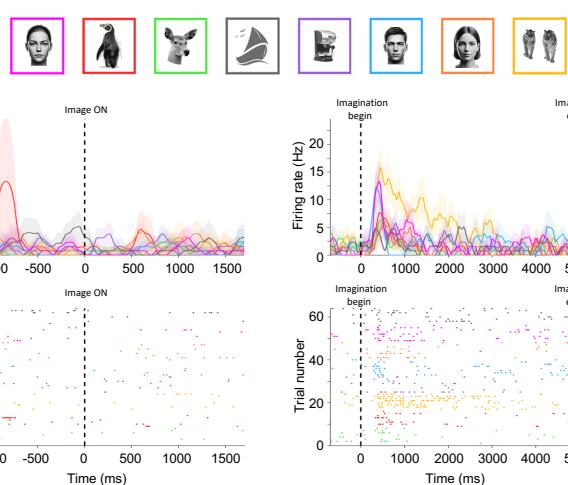
1038 **Figure 3. Decoding object features via linear regression and generating super stimuli for VTC**
1039 **neurons using the axis model and a GAN.**

1040 [Note: The human face in panel H of this figure has been replaced with a synthetic face
1041 generated by a diffusion model (102), in accordance with the bioarxiv policy on displaying
1042 human faces.]

1043 (A-E) Reconstruction of objects from neural responses. **(A)** Decoding model. We used responses
1044 to all but one object ($500-1 = 499$) to determine the transformation between responses and
1045 feature values by linear regression, and then used that transformation to predict the feature
1046 values of the held-out object. **(B)** Decoding accuracy. Predicted vs. actual feature values for the
1047 first (left, $p = 1.04\text{e-}25$, paired t-test) and second (right, $p = 0$, paired t-test) dimensions of object
1048 space (see methods). **(C)** Distribution of normalized distances between reconstructed feature
1049 vectors and the best-possible reconstructed feature vectors for 482/500 images (see methods) to
1050 quantify decoding accuracy across the population. The normalized distance takes into account
1051 the fact that the object images used for reconstruction did not include any of the object images
1052 shown to the patients. A normalized distance of 1 means the reconstruction found the best
1053 solution possible. **(D)** Images were split into tertiles based on normalized distance. Examples of
1054 the reconstructions in the first tertile (top row), second (middle row) and third (bottom row) as
1055 compared to the original stimulus images being reconstructed. **(E)** Decoding accuracy as a
1056 function of the number of distractor objects drawn randomly from the stimulus set (see
1057 methods). The black dashed line represents the decoding accuracy one would expect by chance.
1058 **(F)** Axis plot of an example neuron showing the positions of generated stimuli sampled along the
1059 preferred axis (yellow) and orthogonal axis (green) relative to the other stimulus images. The
1060 maximum and minimum projection valued images are displayed along the preferred axis. The
1061 vertical and horizontal line plots are the binned firing rate of the stimulus images as one moves
1062 along each axis with the generated images overlaid on top. The systematic increase along the
1063 preferred axis with an almost identical response to all images along the orthogonal axis is clearly
1064 visible. **(G)** (left) The responses of the neuron in (F) to a grid of images sampled in the space
1065 spanned by the preferred and orthogonal axes. The extent of the grid is indicated by the purple
1066 bounding box in (F). (right) The stimuli that comprised the grid. The neuron showed increases in
1067 firing rate to stimuli that varied in directions parallel to the preferred axis. **(H)** Another example
1068 neuron with responses to both screening and synthetic stimuli (see (F) for notation). **(I)** Population
1069 summary. For each of the 16 neurons across 4 patients, the Pearson Correlation between the
1070 projection value of the generated images and the firing rate response of the neuron is visualized
1071 for the preferred and orthogonal axes respectively. The two example neurons shown in (F-G) and
1072 (H) are marked in red. **(J)** Distribution of the correlation values between the predicted and
1073 observed firing rate responses to the synthetic images. For each neuron, the preferred axis was
1074 computed using the 500 stimulus images and used to learn the transformation between
1075 projection values and firing rates (as in A, see methods). This transformation was used to compute
1076 predicted responses of the neuron to the synthetic images. The distribution shown is the
1077 correlation of those predicted values (using the axis of the neuron in the first session) to the
1078 responses of the matched neuron in the second session observed during the experiment (see
1079 methods for matching procedure).

A

Total number of neurons = 231

Active during imagery
= 66/231Silent during imagery
= 165/231**B****C****D**

1080 **Figure 4. Cued Imagery task reveals robust reactivation of VTC neurons.**

1081 [Note: The human faces in panels B-D of this figure have been replaced with synthetic faces
1082 generated by a diffusion model (102), in accordance with the bioarxiv policy on displaying
1083 human faces.]

1084 **(A)** The proportion of total recorded VTC neurons that were active (or reactive) during imagery.
1085 **(B-D)** Example neurons. Left PSTH and raster shows the response of the neuron during the
1086 encoding while right PSTH and raster shows the response during imagery. The stimulus images
1087 used in the task are shown above arranged in ascending order of projection onto the neurons
1088 preferred axis (computed during screening). **(B)** This neuron's preferred stimulus was the laptop
1089 during encoding and it reactivated robustly during imagery of the laptop. Note that this was not
1090 an axis tuned neuron and as such the laptop is not the right most image. **(C)** Example of an axis
1091 tuned unit that reactivated to multiple stimuli (red & green) in a graded manner during imagery.
1092 **(D)** A small number of neurons recorded (15/231) were quiet during encoding but strongly active
1093 during imagery. Once such example is shown.

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

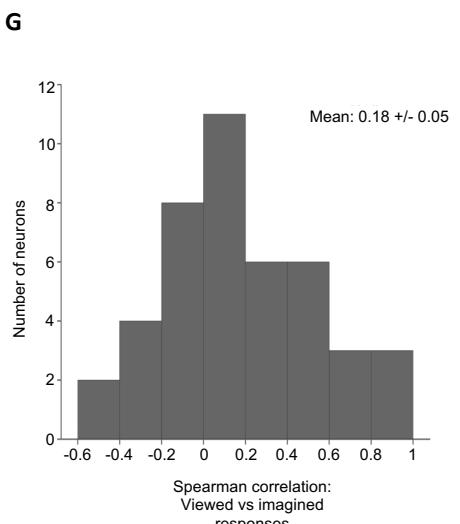
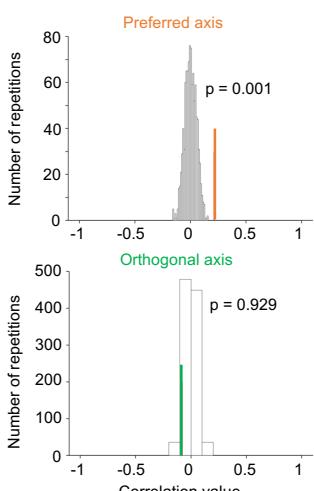
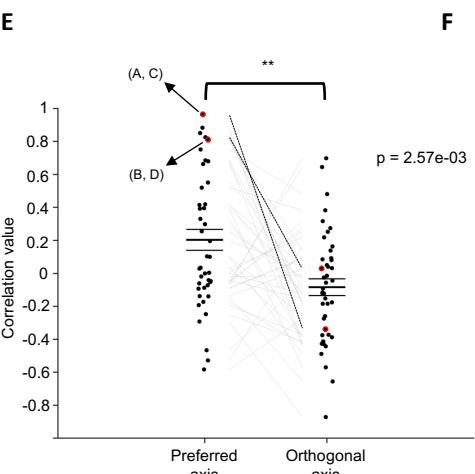
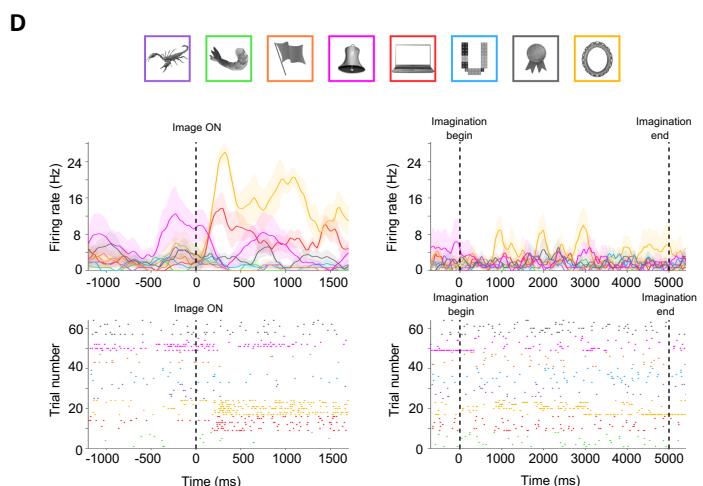
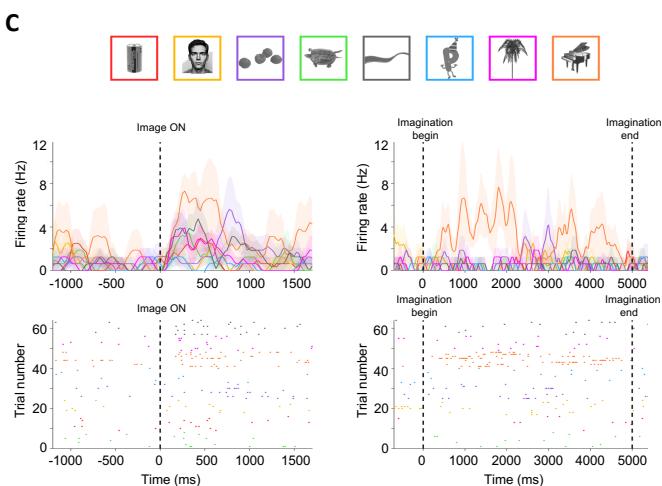
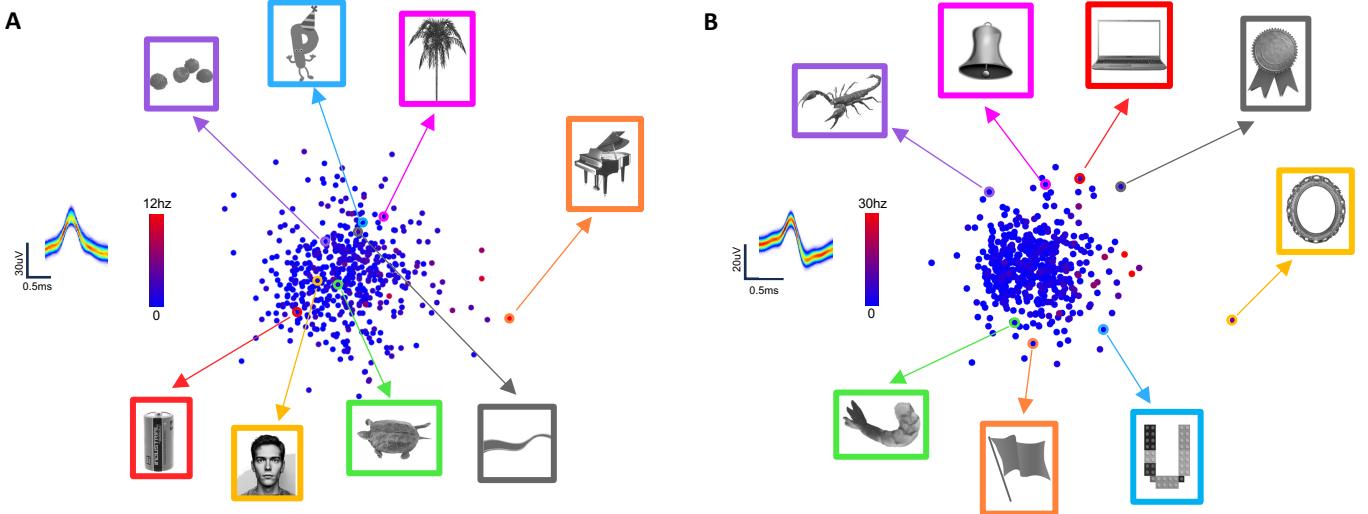
1104

1105

1106

1107

1108



1109 **Figure 5. Neurons in VTC reactivate in a manner that respects the visual code.**

1110 [Note: The human face in panels A&C of this figure have been replaced with a synthetic face
1111 generated by a diffusion model (102), in accordance with the bioarxiv policy on displaying
1112 human faces.]

1113 **(A-D) Two example neurons.** **(A, B)** Axis plots with the subset of 8 stimuli used for imagination
1114 indicated. Inset shows the waveform. **(C, D)** Response during encoding/viewing (left) and
1115 imagery (right). The top panel shows the stimuli, arranged in order of increasing projection
1116 value along the preferred axis. **(E-F)** Population summary. **(E)** Pearson correlation between the
1117 projection value onto the axis computed during screening and firing rate during imagery is
1118 shown for the preferred and orthogonal axes (n=43 neurons). Across neurons the preferred
1119 axes showed a significantly higher positive value ($r_{pref} = 0.20$, $r_{ortho} = -0.084$, $p = 2.57e-03$
1120 Wilcoxon ranksum test). The neurons discussed in (A-B) and (C-D) are marked in red. **(F)**
1121 Comparison of the mean correlation across all reactivated neurons to the null distribution.
1122 Responses were significantly correlated to the projection value onto the preferred axis ($p =$
1123 0.001, shuffled distribution with 1000 repetitions) but not to the projection value onto the
1124 orthogonal axis ($p = 0.929$, shuffled distribution with 1000 repetitions). **(G)** Correlation between
1125 the firing rates during viewing and imagination of the same stimuli. The responses to each
1126 stimulus in encoding and imagery were averaged across trials and the Pearson correlation
1127 coefficient was computed between those two vectors for each neuron. The mean value is 0.18,
1128 significantly larger than 0 ($p = 2.40e-3$, one sample ttest).

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145
1146
1147
1148
1149
1150 Supplementary materials for
1151
1152 **A shared code for perceiving and imagining objects in human ventral**
1153 **temporal cortex**

1154
1155 V. S. Wadia, C. M. Reed, J. M. Chung, L. M. Bateman, A. M. Mamelak,
1156 U. Rutishauser[¶], D. Y. Tsao[¶]

1161 **This PDF file includes:**

1164
1165 Supplementary Results
1166 Supplementary Materials and Methods
1167 Figs. S1 – S6
1168 Table S1, S2

1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179

1180 **Supplementary results**

1181

1182 **Axis code is independent of the specific convolutional network used to**

1183 **parametrize stimuli (relevant for Figures S3 & S4)**

1184

1185 All analyses reported in the main text use a feature space built from layer fc6 of AlexNet.

1186 However, there now exist a plethora of deep convolutional neural network models that achieve

1187 high performance on object recognition (1). We therefore set out to compare the ability of

1188 several such models to explain the responses of VTC neurons to general objects.

1189

1190 The models tested include AlexNet (2), both VGG-16 and 19 (3), VGG-Face (4), the eigen object

1191 model wherein the space is built by performing PCA on the pixel level representation directly (5,

1192 6), and four CORNet models (7, 8). AlexNet, VGG-16/19, and the CORNet family of networks are

1193 trained to classify images into 1000 object categories with varying architectures. AlexNet has 8

1194 layers, with 5 convolutional and 3 fully connected layers. VGG-16 and 19 have 16 and 19 layers

1195 respectively, with 3 fully connected layers and the rest being convolutional layers. VGG models

1196 are also known for leveraging the smallest possible receptive field in their convolutional layers

1197 (3x3). The CORNet family of networks consists of 4 networks: CORNet-Z is purely feedforward;

1198 CORNet-R includes some recurrence which has been shown to be essential for object

1199 recognition in the primate visual system (9); CORNet-RT has the same structure as R but

1200 includes ‘biological unrolling’ wherein the input at time $t + 1$ in layer n is the same as input to

1201 layer $n - 1$ at time t so that information flows through the layers sequentially (7); and finally

1202 CORNet-S has the most complicated of architecture, including recurrent and skip connections

1203 between the layers (8). Despite the individual differences all four networks have architectures

1204 inspired by the primate visual system with layers corresponding to V1, V2, V4, and IT. VGG-Face

1205 has the same architecture as VGG-16 but is trained to identify 2622 celebrities (4).

1206

1207 To quantify the ability of each network to explain human VTC responses we learned a linear

1208 mapping between the features of each model and the neural responses (10). As we did in our

1209 earlier axis tuning computations and to avoid overfitting, we reduced dimensionality of the

1210 feature representations via principal-component analysis (PCA) yielding 50 features for each

1211 object and model. In our main analyses, we used leave-one-out cross validation and for each

1212 neuron fit the responses of 500-1 images to the 50 features via linear regression before

1213 predicting the response of the neuron to the left-out image using the same linear transform.

1214 The explained variance by the linear transform was used as an initial measure of goodness-of-

1215 fit. Beyond this we also computed the encoding and decoding error for each neuron with every

1216 model. Encoding error was computed as follows: the observed population vector to each object

1217 was compared to the predicted population response vector and the observed population

1218 response vector to a random other object in the set. If the angle between the observed and

1219 predicted response to the chosen object was smaller than the angle between the predicted

1220 response and the distractor, the prediction was considered correct (Figure S3B). Decoding error

1221 was computed via the same method except the feature vector of the object was predicted from

1222 the neural responses. In other words, the roles of the neural responses and the model features
1223 were reversed (see methods). A model was considered to explain neural responses well if it has
1224 high explained variance and low encoding/decoding error. We found that with the exception of
1225 VGG-Face and the eigen object model that performed significantly worse, there was no
1226 significant difference in explained variance between the models (Figure S3A; $p = 8.72\text{e-}10$,
1227 AlexNet vs VGG-Face, Wilcoxon ranksum test; $p = 4.49\text{e-}25$, AlexNet vs eigen model, Wilcoxon
1228 ranksum test). The most complicated CORNet, CORNet-S outperformed the purely feedforward
1229 CORNet-Z ($p = 1.69\text{e-}3$, CORNet-S vs CORNet-Z, Wilcoxon ranksum test) but not its recurrent
1230 counterparts ($p = 0.59$, CORNet-S vs CORNet-R; $p = 0.601$, CORNet-S vs CORNet-RT, Wilcoxon
1231 ranksum test).

1232

1233 A note on the heterogeneity of mental representations across people (relevant for
1234 Figure S6)

1235

1236 There is a large body of evidence to support the notion that the subjective vividness of visual
1237 imagery varies greatly between individuals (11–13), with some individuals demonstrating a
1238 complete inability to generate a mental image (aphantasia) (14) while others have near-
1239 photorealistic mental images (hyperphantasia). Moreover, various neuroimaging studies have
1240 shown differences in fMRI bold signals — both intensity in early visual areas (15) and functional
1241 connectivity between areas (12, 16) — between subjects that reported different amounts of
1242 vividness in imagery. Growing evidence of these differences has led to the conclusion that
1243 examining mental imagery at the group level with current tools (fMRI and psychophysics) is not
1244 appropriate — leading to the end of the “imagery debate” (17, 18).

1245

1246 In order to understand whether there is a correlation between the data discussed here and
1247 subjective vividness, patients also completed the Vividness of Visual Imagery Questionnaire
1248 (VVIQ) (19). These responses were recorded in-person starting with the 3rd patient in this study
1249 and retroactively via video call for the previous ones. Remarkably, all the patients discussed
1250 here had very high scores in the vividness scale, with every single one of them falling into the
1251 ‘hyperphantasic’ category (Figure S6D). It therefore remains unclear whether the recapitulation
1252 of sensory context demonstrated in this study extends to individuals with weak visual imagery
1253 capabilities.

1254

1255

1256

1257

1258

1259

1260

1261

1262 **Supplementary Methods**

1263 **Model comparisons**

1265 **Extraction of features from stimulus images**

1268 Each stimulus image was fed into one of the following models to extract the corresponding
1269 features:

1270 **Eigen object model:**

1273 PCA was performed on the original images of the 500 stimulus objects and the top 50 PCs were
1274 extracted to compare with other models.

1275

1276 **AlexNet:**

1278 We used a pre-trained MATLAB implementation of AlexNet. This is an 8 layer deep convolutional
1279 neural network with 5 convolutional layers and 3 fully connected layers, trained to classify
1280 images into 1000 object categories.

1281

1282 **VGG Family:**

1284 We used pre-trained MATLAB implementations of VGG-16, a 16 layer deep convolutional neural
1285 network that contains 16 layers with 13 convolutional layers and 3 fully connected layers trained
1286 to classify images into 1000 object categories (3), VGG-Face which has the same structure as
1287 VGG-16 but is trained to recognize the faces of 2622 celebrities (4), and VGG-19 which has 19
1288 layers (16 convolutional and 3 fully connected) trained on the same task as VGG-16 (3).

1289

1290 **CORNet Family:**

1292 We used a pre-trained PyTorch implementation of CORNet. The CORNet family contains 3
1293 architectures: CORNet-Z, CORNet-R, and CORNet-S. Each architecture includes 4 main layers that
1294 correspond to V1, V2, V4, and VTC. CORNet-Z is the simplest model and is purely feedforward.
1295 CORNet-R takes the otherwise feedforward network and introduces recurrent dynamics within
1296 each area. CORNet-S is the most complex containing within area recurrent connections, skip
1297 connections and the most convolutional layers. Our plots include a CORNet-RT plot which refers
1298 to a version of CORNet-R that does biological temporal unrolling (7) (see fig 2 in reference).

1299
1300 The parameters of AlexNet and VGG were obtained from MATLAB's Deep Learning Toolbox. The
1301 CORNets were downloaded from (<https://github.com/dicarlolab/CORnet>).

1302

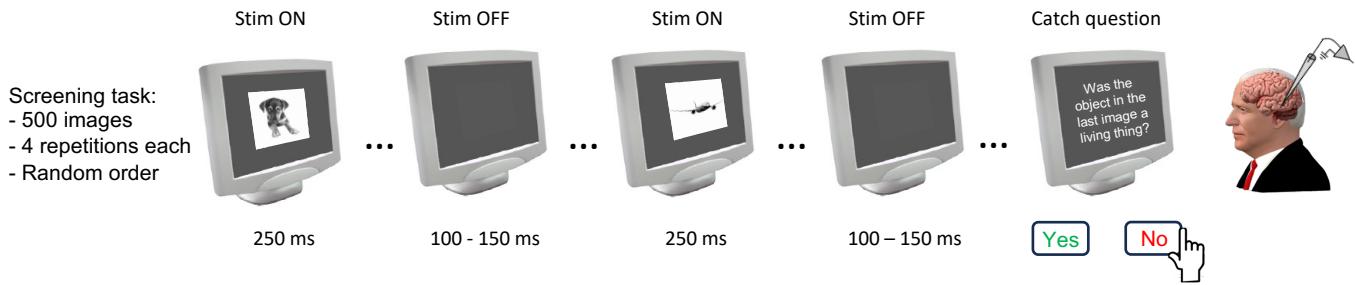
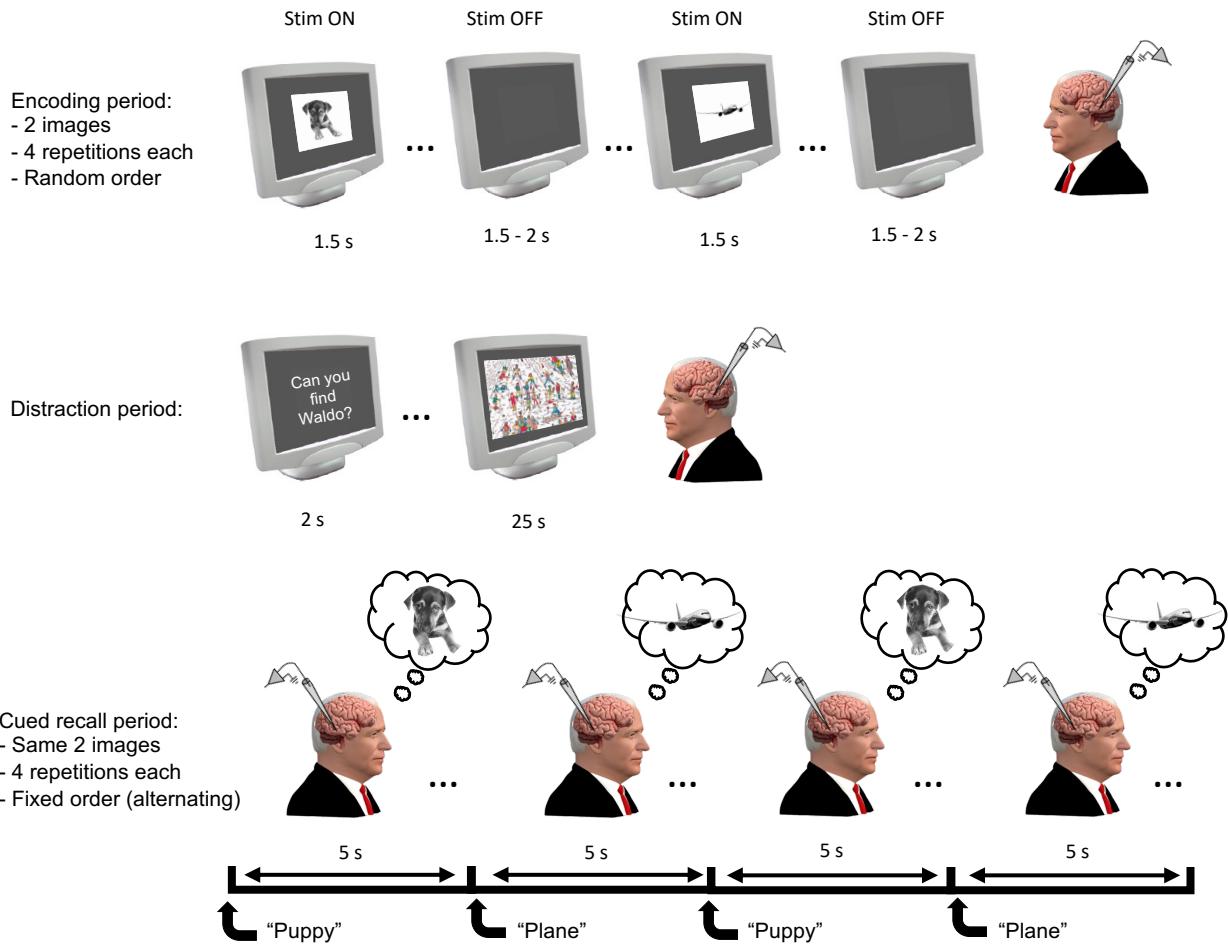
1303

1304 **Explained variance computation**
1305
1306 See main methods ('Axis model' subheading). The reported results are the ratio of explained to
1307 explainable variance in the 321/367 axis tuned neurons that had a > 10% explainable variance.
1308
1309 **Encoding/Decoding error computation**
1310
1311 For the encoding analysis, the response of each neuron was first zscored, then the same
1312 procedure as the explained variance computation was followed to obtain a predicted response
1313 to every single object. To quantify prediction accuracy, we examined the angle between the
1314 predicted population response to each object and its actual population response (target) or the
1315 population response to a different object (distractor). If the angle between the predicted
1316 response vector and the distractor was smaller than the angle between the predicted response
1317 vector and the target this was counted as an error. Overall encoding error was quantified as the
1318 average errors across 1000 pairs of target and distractor objects.
1319
1320 For the decoding analysis we used exactly the same procedure, however the roles of the neural
1321 responses and the object features were reversed. We first normalized each dimension of object
1322 features to have zero mean and unit variance, then for each image we used a leave-one-out
1323 procedure to fit a linear transform using the responses to 499 images and then predict the
1324 features of the left-out image. Decoding error was computed as the average decoding error
1325 across all target and distractor pairs in feature space.
1326
1327
1328 **Quantification of axis consistency**
1329
1330 The consistency of the preferred axis of a neuron (Figure S2E) was determined as follows: The
1331 image set was randomly split into two subsets of 796 and 797 images, a preferred axis was
1332 computed for each set and the Pearson correlation was computed between the two. This
1333 procedure was repeated 100 times and the axis consistency was defined as the to get an
1334 aggregate correlation value.
1335
1336
1337
1338
1339
1340
1341
1342

1343 **Supplementary References**

- 1344
- 1345 1. Y. Taigman, M. Yang, M. Ranzato, L. Wolf, “DeepFace: Closing the gap to human-level
1346 performance in face verification” in *2014 IEEE Conference on Computer Vision and Pattern
1347 Recognition* (IEEE, 2014), pp. 1701–1708.
- 1348 2. A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional
1349 neural networks. *Adv. Neural Inf. Process. Syst.* **25** (2012).
- 1350 3. K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image
1351 Recognition, *arXiv [cs.CV]* (2014). <http://arxiv.org/abs/1409.1556>.
- 1352 4. O. Parkhi, A. Vedaldi, A. Zisserman, Deep face recognition. *BMVC 2015 - Proceedings of the
1353 British Machine Vision Conference 2015* (2015).
- 1354 5. L. Sirovich, M. Kirby, Low-dimensional procedure for the characterization of human faces. *J.
1355 Opt. Soc. Am. A* **4**, 519–524 (1987).
- 1356 6. M. A. Turk, A. P. Pentland, Face Recognition Using Eigenfaces (1991).
1357 <https://www.cin.ufpe.br/~rps/Artigos/Face%20Recognition%20Using%20Eigenfaces.pdf>.
- 1358 7. J. Kubilius, M. Schrimpf, A. Nayebi, D. Bear, D. L. K. Yamins, J. J. DiCarlo, CORnet: Modeling
1359 the Neural Mechanisms of Core Object Recognition, *bioRxiv* (2018)p. 408385.
- 1360 8. J. Kubilius, M. Schrimpf, K. Kar, R. Rajalingham, H. Hong, N. Majaj, E. Issa, P. Bashivan, J.
1361 Prescott-Roy, K. Schmidt, Others, Brain-like object recognition with high-performing
1362 shallow recurrent ANNs. *Adv. Neural Inf. Process. Syst.* **32** (2019).
- 1363 9. K. Kar, J. Kubilius, K. Schmidt, E. B. Issa, J. J. DiCarlo, Evidence that recurrent circuits are
1364 critical to the ventral stream’s execution of core object recognition behavior. *Nat. Neurosci.*
1365 **22**, 974–983 (2019).
- 1366 10. L. Chang, B. Egger, T. Vetter, D. Y. Tsao, Explaining face representation in the primate brain
1367 using different computational models. *Curr. Biol.* **31**, 2785-2795.e4 (2021).
- 1368 11. J. Pearson, T. Naselaris, E. A. Holmes, S. M. Kosslyn, Mental Imagery: Functional
1369 Mechanisms and Clinical Applications. *Trends Cogn. Sci.* **19**, 590–602 (2015).
- 1370 12. N. Dijkstra, S. E. Bosch, M. A. J. van Gerven, Vividness of Visual Imagery Depends on the
1371 Neural Overlap with Perception in Visual Areas. *J. Neurosci.* **37**, 1367–1373 (2017).
- 1372 13. F. Galton, Visualised Numerals, *Nature Publishing Group UK* (1880).
1373 <https://doi.org/10.1038/021252a0>.
- 1374 14. A. Zeman, M. Dewar, S. Della Sala, Lives without imagery - Congenital aphantasia. *Cortex*
1375 **73**, 378–380 (2015).

- 1376 15. X. Cui, C. B. Jeter, D. Yang, P. R. Montague, D. M. Eagleman, Vividness of mental imagery:
1377 individual variability can be measured objectively. *Vision Res.* **47**, 474–478 (2007).
- 1378 16. J. Fulford, F. Milton, D. Salas, A. Smith, A. Simler, C. Winlove, A. Zeman, The neural
1379 correlates of visual imagery vividness – An fMRI study and literature review. *Cortex* **105**,
1380 26–40 (2018).
- 1381 17. S. M. Kosslyn, Image and brain: the resolution of the imagery debate. *J. Cogn. Neurosci.* **7**,
1382 415–420 (1995).
- 1383 18. J. Pearson, S. M. Kosslyn, The heterogeneity of mental representation: Ending the imagery
1384 debate. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 10089–10092 (2015).
- 1385 19. D. F. Marks, Vividness of Visual Imagery Questionnaire. *British Journal of PsychologyJournal*
1386 *of Mental Imagery*, doi: 10.1037/t05959-000 (1973).
- 1387 20. A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, I. Sutskever, Zero-shot
1388 text-to-image generation. *ICML abs/2102.12092*, 8821–8831 (2021).
- 1389 21. P. Bao, L. She, M. McGill, D. Y. Tsao, A map of object space in primate inferotemporal
1390 cortex. *Nature* **583**, 103–108 (2020).
- 1391 22. D. L. K. Yamins, H. Hong, C. F. Cadieu, E. A. Solomon, D. Seibert, J. J. DiCarlo, Performance-
1392 optimized hierarchical models predict neural responses in higher visual cortex. *Proc. Natl.*
1393 *Acad. Sci. U. S. A.* **111**, 8619–8624 (2014).
- 1394 23. C. F. Cadieu, H. Hong, D. L. K. Yamins, N. Pinto, D. Ardila, E. A. Solomon, N. J. Majaj, J. J.
1395 DiCarlo, Deep neural networks rival the representation of primate IT cortex for core visual
1396 object recognition. *PLoS Comput. Biol.* **10**, e1003963 (2014).
- 1397
- 1398
- 1399
- 1400
- 1401
- 1402
- 1403
- 1404
- 1405

A**B**

1406 **Figure S1. Detailed schematics of both screening and cued imagery tasks.**

1407 Outlines showing the task structure for both screening and cued imagery tasks including
1408 number of images used, different task stages, and stimulus order.

1409 **(A)** Screening task schematic. In the screening task, grayscale images with white backgrounds
1410 were displayed on a gray screen for 250 ms with the inter-trial interval jittered between 100-150
1411 ms. Images subtended 6-7 visual degrees. At random intervals (min interval: 1 trial, max
1412 interval: 80 trials) a yes-no catch question would appear pertaining to the image that came just
1413 before it. Each image was repeated 4 times for a total of 2000 trials **(B)** Cued imagery task
1414 schematic. In the cued imagery task, a subset (6-8) of the 500 images used for screening were
1415 used, chosen to have spread across both the preferred and principal orthogonal axes. Each trial
1416 consisted of 2 images, and an initial encoding period wherein images were displayed on a gray
1417 screen for 1.5 s with the inter trial interval jittered between 1.5 – 2 s. After each image was
1418 viewed 4 times, a distraction period occurred wherein patients were required to spend 30 s on a
1419 visual search puzzle. Finally, after the distraction period they were cued verbally by the
1420 experimenter to imagine both images in the trial one by one in alternating order until both had
1421 been visualized 4 times. Each image appeared in 2 trials for a total of 8 repetitions of viewing
1422 and imagery per image.

1423

1424

1425

1426

1427

1428

1429

1430

1431

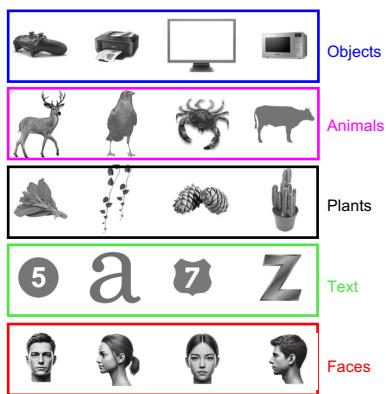
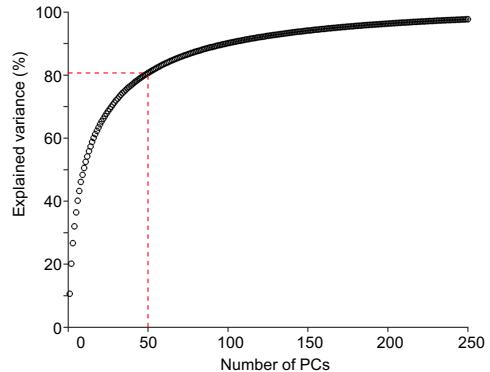
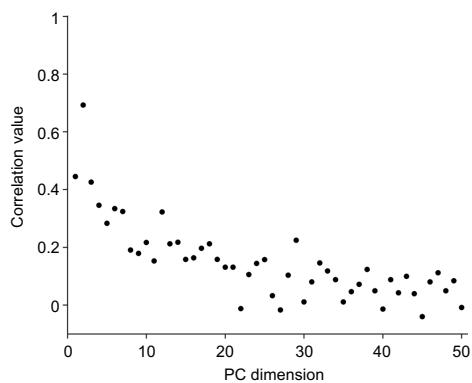
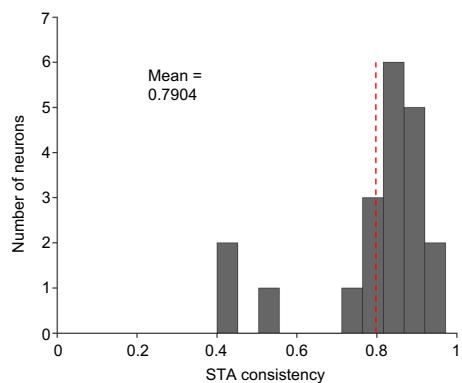
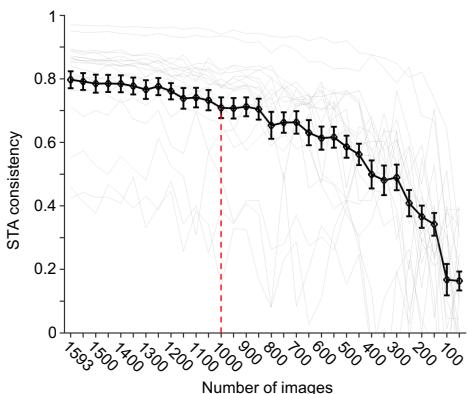
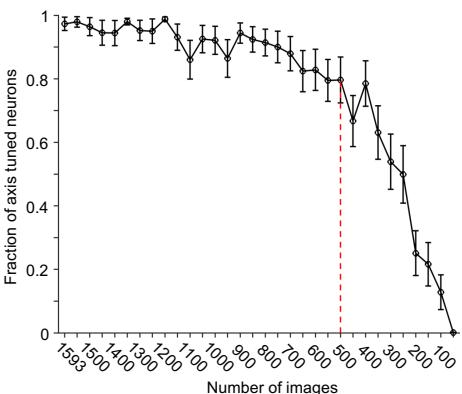
1432

1433

1434

1435

1436

A**B****C****D****E****F**

1437 **Figure S2. Stimuli and parameters chosen for screening task.**

1438 [Note: The human faces in panel A of this figure have been replaced with synthetic faces
1439 generated by a diffusion model (20), in accordance with the bioarxiv policy on displaying human
1440 faces.]

1441 **(A)** 500 stimuli from face, text, plant, animal, and object categories were shown to the patients.
1442 Example images from each of the 5 categories (taken from www.freepngs.com). **(B)** The
1443 cumulative explained variance of fc6 unit responses over 250 PCs. 50 dimensions explained
1444 80.78% of the variance. **(C)** Correlation between predicted feature values using the axes
1445 computed via the 500 stimulus screening and actual feature values for all 50 dimensions (related
1446 to Figure 3B) **(D-F)** In one early session a patient performed a more comprehensive version of the
1447 screening task, with 1593 stimuli (21). The 500 used in this project were subsampled from this
1448 larger set. The responses of the 22 axis tuned neurons in this session were used to determine
1449 appropriate parameters going forward. **(D)** Distribution of axis consistency (computed using half-
1450 splits, see methods) for axis tuned neurons in this session. Red line indicates the mean value. **(E)**
1451 Axis consistency as a function of the number of stimuli used to compute axes. The red line at 1000
1452 indicates how consistent axes would be if computed using 500 stimuli. Axis consistency remains
1453 stable until the stimulus number drops below 400 in each half split. This informed the choice of
1454 500 stimuli. **(F)** Proportion of axis tuned neurons detected as a function of the number of stimuli
1455 used. 78% of the axis tuned neurons detected using 1593 stimuli were still detected using 500
1456 stimuli. At lower stimulus numbers the neuron count drops precipitously.

1457

1458

1459

1460

1461

1462

1463

1464

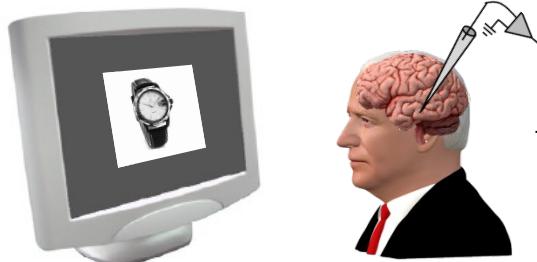
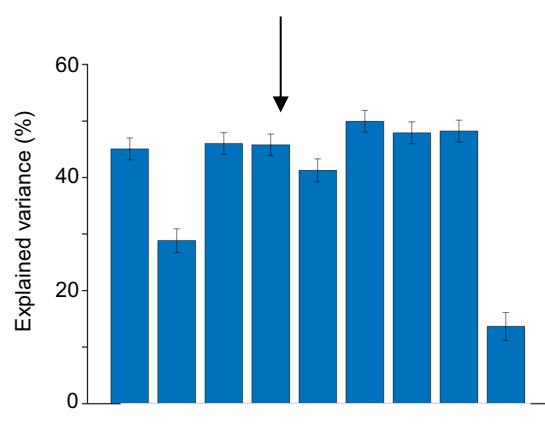
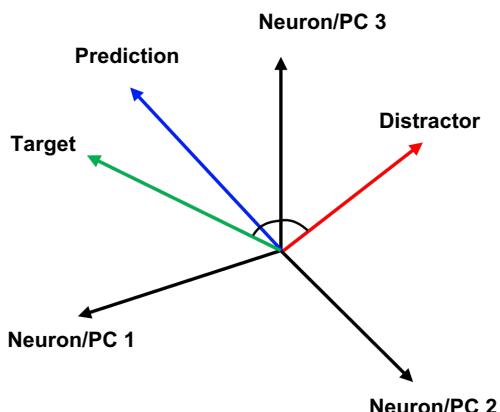
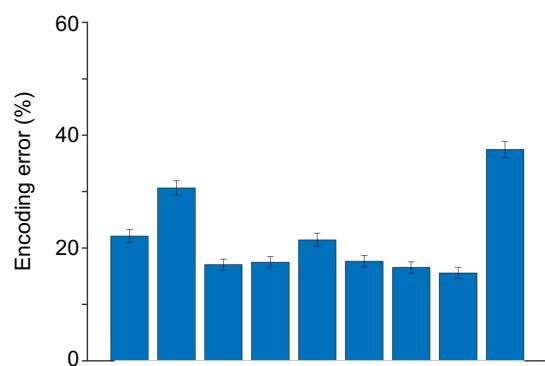
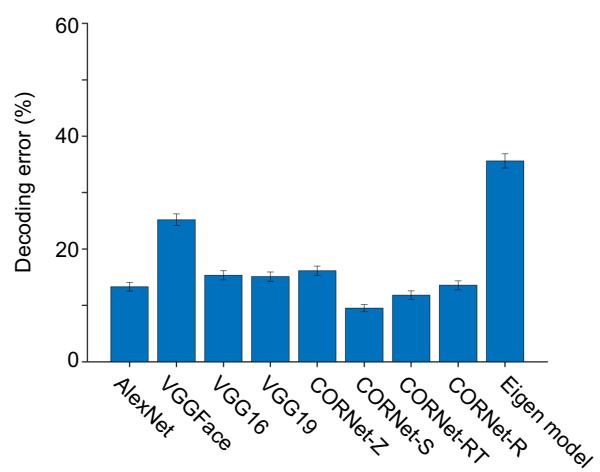
1465

1466

1467

A

- AlexNet (Krizhevsky et al., 2012)
 Eigen Object Model (Sirovich & Kirby, 1987)
 VGG-16 (Simonyan & Zisserman, 2015)
 VGG-19 (Simonyan & Zisserman, 2015)
 VGG-Face (Parkhi et al., 2015)
 CORNet-Z (Kubilius et al., 2018)
 CORNet-R/RT (Kubilius et al., 2018)
 CORNet-S (Kubilius et al., 2019)

B**C****D****E****F**

1468 **Figure S3. Comparing how well various models explained human VTC responses to object**
1469 **images.**

1470 **(A)** The 500 grayscale stimulus images used in the screening task were parametrized using 8
1471 different models from 4 different model families (AlexNet, the eigen object model, VGG, and
1472 CORNet). The same number of features were extracted from units of the different models using
1473 principal components analysis (PCA) for comparison. **(B)** Responses from VTC neurons were
1474 recorded as patients viewed these objects 4 times each. For recording locations and task
1475 schematic see Figure 1D and Figure S1A respectively. **(C)** The explained variances for each model
1476 after 50 features were extracted using PCA. For each neuron, explained variance was normalized
1477 by the explainable variance (see methods). Error bars represent SEM for the recorded neurons.
1478 The eigen model, VGG-Face, and CORNet-Z performed worse than the other models with no
1479 significant differences between the rest ($p = 8.72e-10$, AlexNet vs VGG-Face, Wilcoxon ranksum
1480 test; $p = 4.49e-25$, AlexNet vs eigen model, Wilcoxon ranksum test). **(D-F)** The various models
1481 were compared with respect to how well they could predict the neuronal responses or the object
1482 features. In both cases a leave-one-out procedure was used to learn and test the transformation
1483 between responses and features. To quantify encoding error for example, for each object we
1484 compared predicted responses to individual objects in the neural state space to the actual
1485 responses to that object and a distractor object. **(D)** If the angle between the predicted response
1486 and the actual response was smaller than the angle between the predicted response and the
1487 distractor the encoding was considered correct. To quantify decoding error, we reversed the roles
1488 of the neural responses and the object features and decoded object features before comparing
1489 the decoded features to the actual features for a given object and a distractor object. **(E)** Encoding
1490 error across all models. **(F)** Decoding error across all models. The eigen model, VGG-Face, and the
1491 purely feedforward CORNet-Z had larger encoding/decoding errors than the other models,
1492 consistent with them explaining much less variance as well.

1493

1494

1495

1496

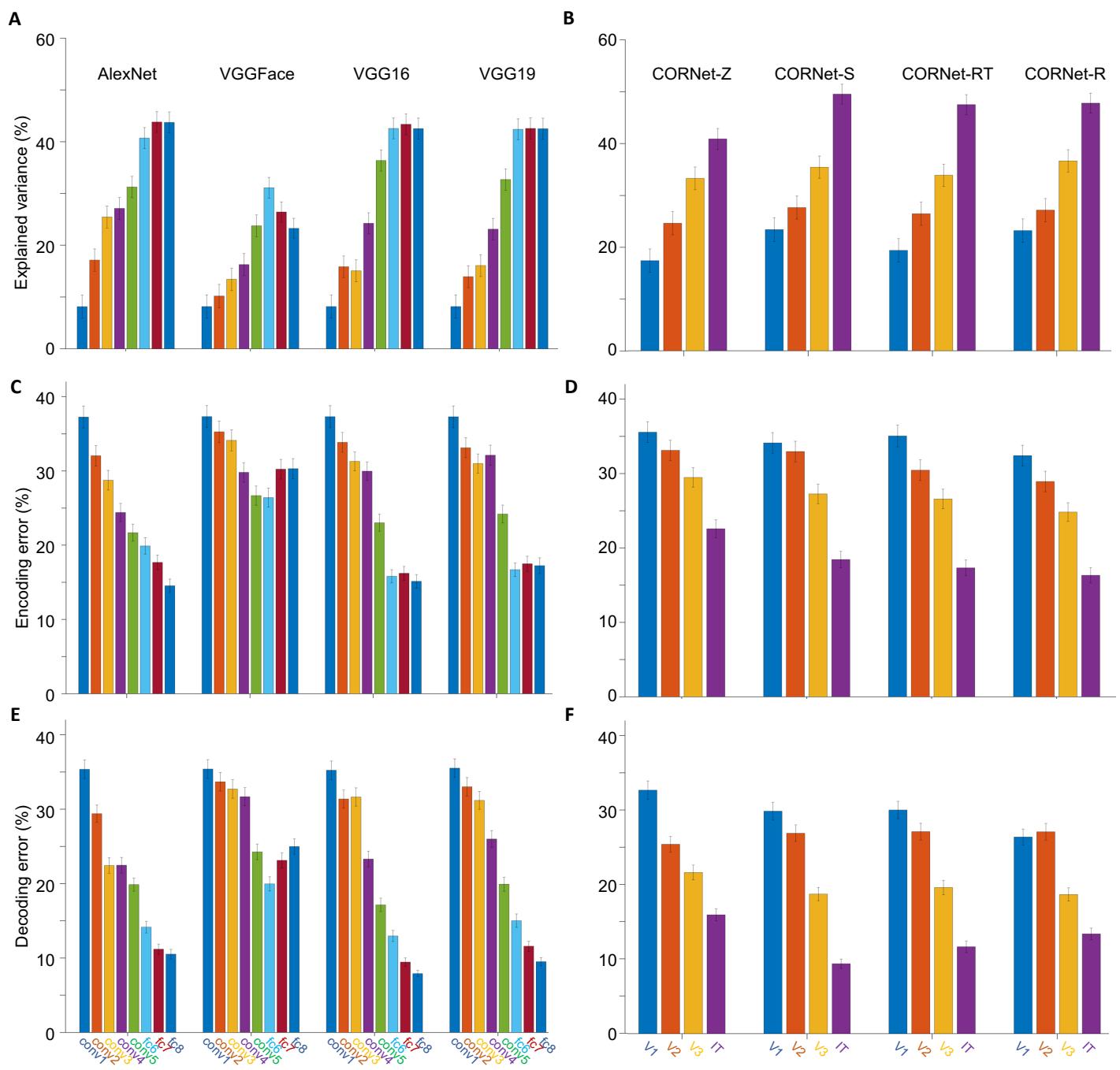
1497

1498

1499

1500

1501



1502 **Figure S4. Comparing explained variance, encoding, and decoding error across layers of**
1503 **various models.**

1504 Inspired by previous work (21–23) units in the penultimate layer were used to build the object
1505 space used in all analyses. Here we chose to compare performance across all fully-connected
1506 layers of AlexNet and the VGG networks, and all output layers of the CORNet networks used in
1507 Figure S3.

1508 **(A)** Explained variance for all the layers of Alexnet and VGG models (before relu, dropout, or
1509 pooling). The performance between the full connected layers is similar. **(B)** Explained variance
1510 across the layers of the various CORNet models (V1, V2, V4, IT). The ‘IT’ or VTC layer performs
1511 the best in all CORNet versions, with little difference across the various forms with the exception
1512 of the purely feedforward version (CORNet-Z) performing worse than its recurrent counterparts.
1513 **(C)** Encoding error across layers of Alexnet and VGG models. As expected, the encoding error is
1514 lowest for the fully connected layers with the highest explained variance. **(D)** Encoding error
1515 across layers of the CORNets. **(E)** Decoding error across layers of Alexnet and VGG networks. **(F)**
1516 Decoding error across layers of CORNets.

1517

1518

1519

1520

1521

1522

1523

1524

1525

1526

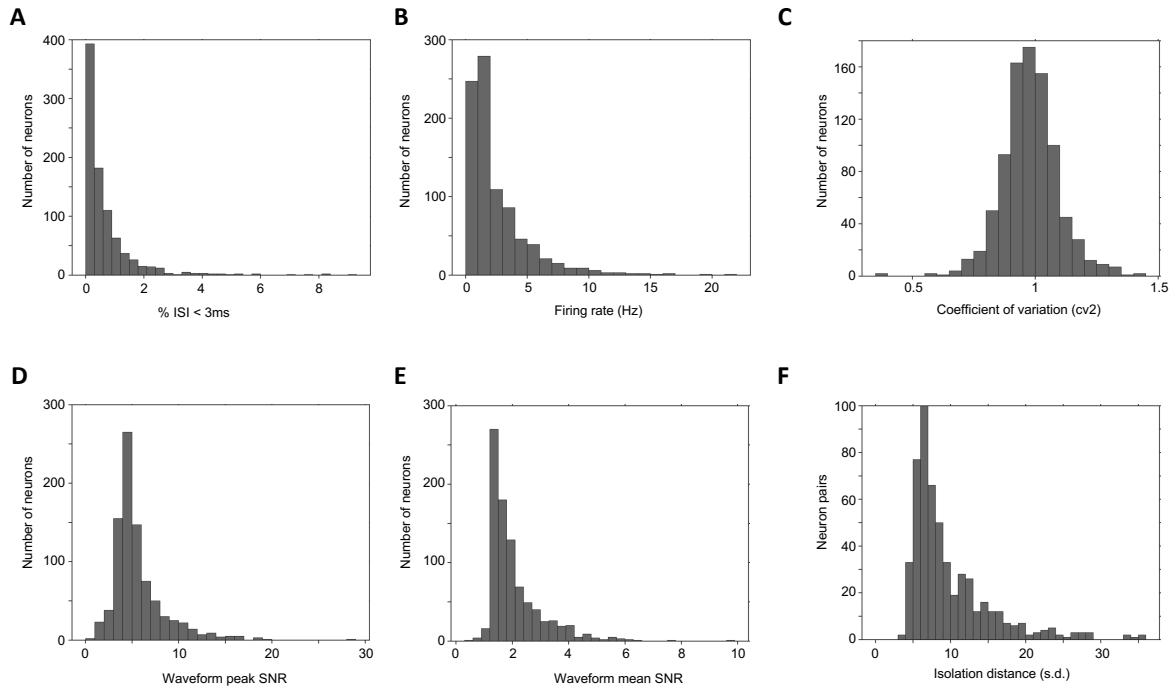
1527

1528

1529

1530

1531



1532 **Figure S5. Spike quality metrics for all identified putative single units.**

1533 **(A)** Proportion of inter-spike intervals (ISI) below 3ms. **(B)** Average firing rate. **(C)** Coefficient-of-
1534 variation. **(D)** Signal-to-noise ratio (SNR) for the peak of the mean waveform across all spikes as
1535 compared to the standard deviation of the background noise. **(E)** Mean SNR of the waveform. **(F)**
1536 Pairwise distance between all pairs of neurons on channels where more than one neuron was
1537 isolated.

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552

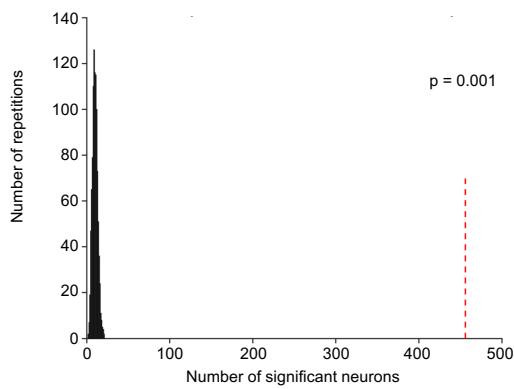
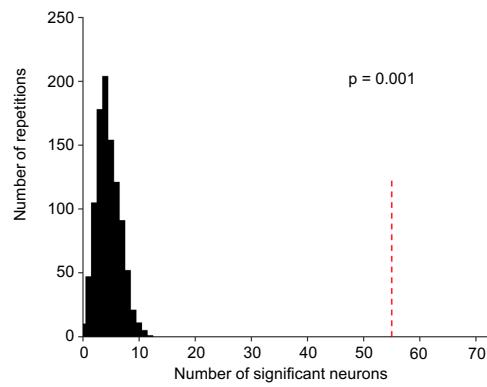
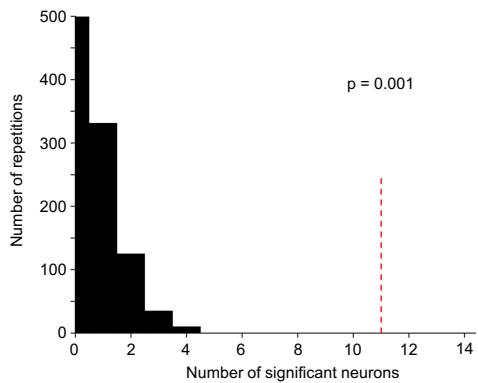
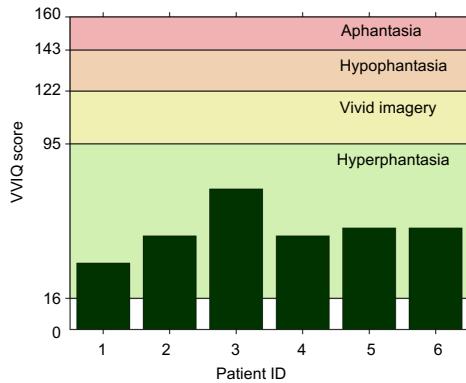
1553

1554

1555

1556

1557

A**B****C****D**

1558 **Figure S6. Bootstrap statistics for number of neurons selected as visually responsive and**
1559 **active during imagery.**

1560 **(A)** Significance of the number of visually responsive VTC neurons. The dashed red line indicates
1561 the number of neurons selected as visually responsive (456/714). The null distribution (black)
1562 was estimated by re-running the identical selection procedure (sliding window ANOVA, see
1563 methods) after randomly shuffling the trial labels. Shuffling the trial labels destroys the
1564 association between the spiking response and trial identity, but keeps everything else (trial
1565 number, stim ON time etc.) intact. The shuffling procedure was carried out for 1000 repetitions.
1566 The mean of the null distribution was 16, implying that the chance level for selecting a neuron
1567 as visually responsive is 16/714 or ~2%. The p-value reported is the percentage of null
1568 distribution values that are greater than the chosen number of neurons. In this case, $p = 0$ and is
1569 reported as 1/number of repetitions $p = 0.001$. **(B, C)** Significance of the number of neurons in
1570 VTC active during imagery. Given that the selection criteria for activation during imagery is
1571 either a number of consecutive significant bins of a sliding window ANOVA or sliding window
1572 ttest we computed a null distribution for each. The null distribution for each is computed by re-
1573 running the identical selection procedure for 1000 repetitions after randomly shuffling the spike
1574 times for the ttest and the trial labels for the ANOVA in the visual responsivity test (see A). The
1575 mean of the null distribution for the ttest (B) was ~5, implying that the chance level for labeling
1576 a neuron as active during imagery via ttest is 5/231 or ~2% and $p = 0.001$, while for the ANOVA
1577 the mean is ~1 so the chance level is 1/231 or ~0.5% and $p = 0.001$. **(D)** All patients completed
1578 the 'Vividness of visual imagery' questionnaire which is a self-assessment of one's visualization
1579 capabilities. All patients recorded from in this study were 'hyperphantasic' or having very vivid
1580 visualizations.

1581

1582

1583

1584

1585

1586

1587

1588

1589

1590

Patient IDs	LVTC (x, y, z)	RVTC (x, y, z)
P73CS	N/A	36.06,-38.29,-13.78
P75CS	N/A	34.12, -55.84, -11.08
P76CS	N/A	34.96, -43.33, -24.04
P77CS	-41.26, -46.26, -15.39	34.03, -40.86, -20.00
P78CS	N/A	28.78, -38.95, -5.94
P79CS	-31.19, -28.97, -24.28	47.59, -29.77, -24.81
P80CS	-39.73, -40.46, -16.97	37.15, -49.70, -12.93
P81CS	-33.27, -61.10, -13.95	29.14, -61.91, -7.52
P82CS	-36.95, -40.17, -17.92	40.58, -42.64, -17.30
P84CS	-31.36, -46.56, -16.80	35.91, -31.69, -25.32
P85CS	-39.76, -41.34, -24.80	N/A
P86CS	-33.16, -40.89, -22.33	28.86, -77.30, -10.03
P87CS	-44.08, -51.67, -23.83	36.87, -48.95, -20.19
P88CS	-43.60, -45.15, -18.00	39.66, -47.79, -15.25
P92CS	-34.38, -41.74, -23.60	38.30, -34.42, -20.74
P93CS	-41.03, -34.84, -27.86	39.40, -38.25, -25.13

1591 **Table S1. MNI coordinates of microwire bundles in which at least 1 VTC neuron was recorded.**
1592 **Related to Figure 1.**

1593

1594

1595

1596

1597

1598

1599

1600

1601

1602

1603

1604

1605

1606

1607

1608

1609

1610

1611

1612

1613

1614

1615

Patient IDs	Sex	Age	Screening		Cued Imagery	
			# Sessions	# resp neurons	# Sessions	# resp neurons
P73CS	F	61	1	16	N/A	N/A
P75CS	F	29	1	6	N/A	N/A
P76CS	F	26	6	22	3	12
P77CS	F	46	1	6	N/A	N/A
P78CS	F	54	2	1	N/A	N/A
P79CS	F	44	7	80	3	44
P80CS	F	25	4	13	2	6
P81CS	F	28	4	38	N/A	N/A
P82CS	M	43	3	3	N/A	N/A
P84CS	M	60	4	47	2	27
P85CS	F	66	6	27	2	9
P86CS	F	39	2	18	N/A	N/A
P87CS	F	27	1	11	N/A	N/A
P88CS	F	29	1	15	N/A	N/A
P92CS	F	30	6	52	2	33
P93CS	F	51	8	101	N/A	N/A

1616 **Table S2. Number of sessions and neurons recorded.**

1617 Summary of the number of neurons recorded in each subject. In some subjects both Screening
1618 and Cued Imagery tasks were performed.

1619

1620