

1
2 **Supplementary Information for**
3 **Texture-like representations of objects in human visual cortex**

4
5
6 Akshay V. Jagadeesh & Justin L. Gardner

7
8 Corresponding Author: Akshay V. Jagadeesh
9 Email: akshayj@stanford.edu

10
11
12 **This PDF file includes:**

13
14 Supplementary text (extended methods)
15 Figure S1: Control analyses using texture stimuli
16 Figure S2: SVM classification accuracy in visual cortex and dCNNs
17 Figure S3: Behavior comparison in-lab vs online
18 Figure S4: Behavior comparison with vs without feedback
19 Figure S5: Pairwise dissimilarity task behavior
20 Figure S6: Natural image selectivity of dCNNs
21 Figure S7: All stimuli used in neuroimaging experiment
22 Figure S8: Synthesized images at different iterations

23
24 References for supplementary text
25
26
27

28

29

30

31

32

33

34 **Extended methods**

35

36 **Stimulus generation**

37

38 All stimuli used in these experiments were either natural images, drawn from image
39 searches of Creative Commons licensed texture and object images, or were synthetically
40 generated through an iterative optimization procedure (“synths”). We selected 34 natural images
41 (22 objects and 12 textures), cropped them into squares, and downsampled the images to 256 x
42 256 pixels. Images were selected as “object stimuli” only if they contained exactly one object,
43 either animate or inanimate, that was clearly visible in the image. Images were selected as
44 texture stimuli if they contained repeated patterns and/or numerous objects that were similar in
45 appearance. Images were categorized as textures or objects by the authors prior to data
46 collection, and these category judgments were confirmed using a Mechanical Turk experiment
47 where an independent sample of 65 subjects were asked to categorize each image as either a
48 texture or object (mean correlation between subjects’ categorizations and our categorization was
49 0.924).

50

51 The image generation procedure, adapted from (1) involved three major steps: feature
52 extraction, spatial pooling, and image synthesis via pixel-wise optimization. In the feature
53 extraction stage, each natural image was passed into an Imagenet-trained VGG-19 deep
54 convolutional neural network (2), and the activations of 3 intermediate layers (pool1, pool2, and
55 pool4) were extracted (**Fig. 1A**). The spatial pooling stage of the standard Gatys algorithm was
56 done by computing the Gramian matrix, i.e. the inner product between all pairs of activation
57 maps, in each layer. The Gramian matrix preserves information about the incidence of individual
58 features as well as the co-incidence of multiple features, while discarding information about the
59 spatial position of those features. Finally, using gradient descent with the L-BFGS algorithmic
60 solver (3), we updated the pixels of a random white noise image to minimize the mean squared
61 error to the Gramian computed for the natural image (1, 4, 5).

61

62 This image synthesis algorithm allowed us to control the complexity of the features in the
63 natural image which are matched in the synthesized output. By varying which layers were
64 included in the loss function, we controlled the complexity of the features in the generated image.
65 This is based on prior research suggesting that early layers of dCNNs encode simple features,
66 such as orientation and spatial frequency, whereas later layers encode more complex features,
67 such as texture, shape, or category identity (6–10). This is an improvement over other texture
68 synthesis algorithms, such as the Portilla-Simoncelli algorithm (4), which includes only one level
69 of higher order statistics computed from the pairwise correlations of a V1-like filterbank (11). We
70 selected 3 different layers from the VGG19 model to include in the synthesis procedure: pool1, an
71 early layer (64 filters); pool2, an intermediate layer (128 filters); and pool4, a late layer (256
72 filters). Layers were added incrementally, so images generated in the pool1 condition include only
73 pool1 features, whereas images generated in the pool4 condition include features from layers
74 pool1, pool2, and pool4.

74

75 This image synthesis algorithm also allowed control over the spatial scale within which
76 the spatial arrangement of features is constrained. Whereas the original Gatys algorithm (1) pools
77 features across the entire image, we modified the algorithm to compute spatially weighted
78 Gramians, which only pooled features within pre-defined spatial pooling regions (4). We tiled the
79 image with equal-sized square spatial pooling regions with smooth transition boundaries defined
80 by a squared cosine function with 20 pixel ramping boundaries. For each unit in the model, we
81 calculated the overlap between its receptive field, calculated based on the kernel size and RF of
82 VGG19, and each spatial pooling region. We used this to compute a spatially weighted Gramian
83 matrix for each pooling region, wherein units are included in proportion to how much their
84 receptive field overlaps the spatial pooling region. By varying the size and number of spatial
85 pooling regions, we imposed stronger or weaker constraints on the spatial arrangement of
86 features. Low spatial constraint synths are ones in which the spatial arrangement of the features
87 can be scrambled across the entire image (which we call 1x1 as there is a single spatial pooling
88 region) and high spatial constraint synths are those in which the arrangement of the features are

88 constrained within small subregions of the image, (for example, a 4x4 set of spatial pooling
89 regions constrains features in subregions that are 1/16th the area of the full image).
90 The number of parameters that constrained each image was a function of the size of the
91 Gramian and the number of spatial pooling regions. The size of the Gramian for a particular layer
92 is equal to the square of the number of filters in each layer, so pool1 images were constrained by
93 4096 (64^2) parameters, pool2 images were constrained by 20480 parameters ($128^2 + 64^2$), and
94 pool4 images were constrained by 282624 parameters ($512^2 + 128^2 + 64^2$). To constrain spatial
95 arrangement of features, we computed a separate Gramian for each spatial pooling region, so the
96 number of spatial pooling regions was a multiplier on the number of parameters. For example, the
97 4x4 images contained 16x the number of parameters as the 1x1 images.

98 Finally, by initiating the optimization process with different random seed images, we
99 generated multiple different synthesis samples which differed significantly in the pixel
100 representation space but contained nearly identical features within each spatial pooling region.
101 For each natural image, we synthesized 3 samples at each of 3 layers (pool1, pool2, and pool4)
102 and 4 spatial constraints (1x1, 2x2, 3x3, 4x4), for a total of 36 synthesized samples per natural
103 image. We used the Adam optimizer (12), implemented in Tensorflow (13), and terminated the
104 image synthesis optimization after 10,000 iterations.

106 **Behavioral Methods**

107 ***Experimental Design – Natural-vs-synth oddity detection task***

110 In the oddity detection experiment, observers performed a 3 alternative forced choice
111 judgment of the odd-one-out (14). On each trial (**Fig. 2A**), observers were asked to fixate
112 centrally on a cross for the duration of the trial, although we could not enforce fixation with eye-
113 tracking and did not employ a central task at fixation. After 200ms of fixation, 3 images were
114 presented -- 1 natural and 2 synths -- concurrently for 2 seconds. Observers were instructed to
115 respond within 2 seconds of stimulus onset, using a keypress, to indicate which image was most
116 different from the others, and on 89.04% of trials, subjects did respond within the time limit (mean
117 RT: 1.08s, SD=0.41s). The two synths were always generated to match the features of the
118 natural image and were both generated from the same layer and spatial constraint, with a
119 different random seed. Following the subject's response, they were shown feedback in the form
120 of the fixation cross changing color for 200ms to either green, indicating a correct response, or
121 red, indicating an incorrect or no response. We also conducted a control experiment where no
122 feedback was given (**Fig. S4**), to ensure that the feedback was not biasing subjects' responses.
123 On each trial, we randomly selected a natural image and 2 different synthesized samples, both
124 with the same feature complexity and spatial constraints, to display. Images subtended
125 approximately 8 degrees, though there was some variability due to the screen and window size of
126 individual participants. Each image was centered 6 degrees away from the fixation cross. We
127 performed this experiment both on Amazon Mechanical Turk, where we recruited 87 subjects
128 who performed a total of 6165 trials, as well as in the lab, where we recruited 2 subjects to
129 perform a total of approximately 5000 trials each and were able to enforce fixation using an
130 Eyelink eyetracking system that aborted any trials where subjects' eye-gaze deviated more than 1
131 degree from the fixation cross. A comparison of in-lab and online data is presented in Supp. Fig.
132 S3. We presented 34 different image classes in this behavioral experiment, including 22 object
133 image classes (**Fig. 2**) and 12 texture image classes (**Fig. S1**), where an image class is defined
134 as the set of images including a natural image and all corresponding feature-matched
135 synthesized samples.

136 ***Experimental Design – Category oddity detection task***

138 To determine human performance at discriminating natural objects of different categories,
139 we recruited human observers to perform a category-level oddity detection task. Like the natural-
140 vs-synth oddity task, subjects were first asked to fixate centrally on a cross for 200ms and were
141 then presented with 3 images concurrently for 2 seconds. Two of those images contained objects
142 from the same category and the third image contained an object of a different category. Subjects

143 were instructed to choose the odd-one-out, i.e. the image which appeared most different from the
 144 others. Images subtended approximately 8 degrees and were centered approximately 6 degrees
 145 away from the fixation cross. We performed this experiment on Amazon Mechanical Turk, where
 146 we recruited 85 subjects who performed a total of 3448 trials.

147 **Experimental Design - Pairwise dissimilarity judgment task**

149 To determine the perceptual similarity of synths with naturals, we conducted a
 150 dissimilarity judgment experiment with an independent set of 110 observers. On each trial,
 151 observers were shown 4 images, grouped into two pairs and were asked to indicate with a
 152 keypress which of the two pairs was more dissimilar. Images subtended 8 degrees. Each pair
 153 was centered 8 degrees to the left and right of fixation, with 4 degrees of vertical separation
 154 between each image. Subjects fixated for 200ms and then stimuli were presented for 2 seconds
 155 and subjects were allowed to respond any time before the images disappeared. No feedback was
 156 given. As with the oddity detection task, all images presented on a given trial were generated to
 157 match the same natural image, and all synths were of the same feature complexity and spatial
 158 constraint. However, unlike the oddity detection task, on a randomly interleaved half of all trials,
 159 all 4 images were synths, and on the other half of the trials, 1 image was the natural image and
 160 the other 3 were synthesized images with scrambled arrangements of features. This enabled us
 161 to determine the perceptual similarity between the synths and the naturals as well as the
 162 perceptual similarity between different synths. We average together all the distances between
 163 pairs of synths to yield a single synth-synth distance. Across all trials, subjects saw 1666 unique
 164 images: 34 image classes \times (1 natural image + (4 synthesized images \times 3 levels of feature
 165 complexity \times 4 levels of spatial constraint)). However, only trials from the 1x1 pool4 condition
 166 were used for estimating perceptual distances. We collected a total of 8687 trials across 110
 167 observers.

168 **Estimating perceptual distances**

170 On any given trial, the observer saw 4 images grouped into 2 pairs, (i_1, i_2) and (i_3, i_4) ,
 171 and was asked to report which pair was more dissimilar. We can thus represent the probability
 172 that the observer will select the first pair (i_1, i_2) as:

$$173 P(D_{1,2} - D_{3,4} + \epsilon > 0)$$

174 where $D_{1,2}$ represents the perceptual distance between the first pair of images, $D_{3,4}$ represents
 175 the distance between the second pair of images, and ϵ is a Gaussian-distributed random variable
 176 with mean 0 and standard deviation σ representing the combination of sensory and response
 177 noise. Then, the probability that the observer will select the first pair is given by $P(\epsilon < D_{1,2} - D_{3,4})$,
 178 which can be computed as the cumulative distribution function of ϵ , $\Phi(x)$ evaluated at $D_{1,2} - D_{3,4}$.
 179 The probability of selecting the second pair is then given by $1 - \Phi(D_{1,2} - D_{3,4})$. Over N trials, if we
 180 observe responses r_1, \dots, r_N , we can compute the likelihood of observing these responses given
 181 the pairwise distances, as

$$182 P(r_1, \dots, r_N | D_{1,2}, D_{1,3}, \dots, D_{N-1,N}, \sigma) = \prod_{i=1}^N \Phi(D_{i_1, i_2} - D_{i_3, i_4})^{r_i} \times (1 - \Phi(D_{i_1, i_2} - D_{i_3, i_4}))^{1-r_i}$$

183 Then, we used the Nelder-Mead optimization algorithm, as implemented in the Python scipy
 184 library (15), to find the values of the distances and the σ that maximize this likelihood function.
 185 For each of the 34 image classes (which were also presented in the oddity task), we estimated
 186 the pairwise distances between 5 images (1 natural, 4 synth), resulting in 10 pairwise distances
 187 $\binom{5}{2}$ to estimate for each image class, yielding a total of 341 parameters (including σ) that were
 188 estimated on 8687 trials of oddity detection behavior.

189 **dCNN observer model**

191
 192 On each trial, our model extracted a feature vector from the last convolutional layer of the
 193 dCNN for each image presented (**Fig. 2B**). Next, we computed the Pearson distance between the

194 features of each pair of images, and for each image, calculated its dissimilarity as the mean
195 Pearson distance from the other two images. Finally, the model converted these dissimilarities
196 into choice probabilities using a Softmax transform. Thus, the probability of choosing the i^{th} item is
197 given by:

198
$$P(c_i) = \frac{e^{\beta c_i}}{\sum_{j=1}^3 e^{\beta c_j}},$$

199 where β is the only estimated parameter, shared across all trials, image classes, and subjects,
200 that is fit to maximize the likelihood of the observed choices and c_i is the mean distance of the i^{th}
201 image from the other two. The β parameter controls the extent to which the model maximizes the
202 choice probability of the most dissimilar image, where a β of 0 yields equal choice probabilities for
203 all images and a beta of infinity would result in a choice probability of 1 for the most dissimilar
204 image. We then visualized these trial-by-trial choice probabilities by computing the average
205 across all the trials of a single condition, to compare the behavior of the model to that of the
206 human subjects.

207 **Modeling IT neurons**

209 To assess the selectivity of neurons in inferior temporal (IT) cortex for natural feature
210 arrangement, we fit a model to a published dataset (16) of multielectrode array recordings
211 measured while macaques passively viewed images of various objects serially presented at the
212 center of gaze. We estimated the response of each neuron as a linear function of activations from
213 each layer of an Imagenet-trained deep convolutional neural network (17, 18), estimated using
214 partial least squares regression, a well-validated approach which yields state-of-the-art
215 predictions of IT neural responses (19, 20). By finding the optimal weighting of dCNN features for
216 best predicting each IT neuron's response, we could then compute a prediction of how each IT
217 neuron would respond to novel images. Then, using this population of 168 model IT neurons, we
218 computed the Pearson distance between the model population's response to each natural image
219 and a corresponding synthesized image as well as the Pearson distance between the model
220 population's response to two different synthesized images of the same class (**Fig. 5A**). Using
221 these two distance measures, we were able to compute a normalized index of selectivity for
222 natural feature arrangement by the formula:

223
$$\frac{d_{\text{natural,synth}} - d_{\text{synth1,synth2}}}{d_{\text{natural,synth}} + d_{\text{synth1,synth2}}}$$

224 Given that our IT model explains, on average, 51.8% of the cross-validated variance in IT
225 neural responses to naturalistic images (18), we cannot treat this as a perfect approximation of IT
226 neurons, although we can use this as a reasonable proxy for IT single unit responses, to
227 corroborate our BOLD imaging evidence. (See Discussion for further consideration of the caveats
228 of this modeling approach).

229 **Neuroimaging Methods**

230 **BOLD Imaging Data Collection**

232 To measure neural responses to natural and synthesized images, we conducted an
233 experiment using blood-oxygen level dependent (BOLD) imaging (21). We recruited seven
234 subjects and instructed them to fixate while visual stimuli were presented over the course of two
235 sessions. To identify the retinotopic maps in visual cortex (22, 23), we presented subjects with
236 four 4-minute runs of a high-contrast sweeping bar stimulus (33), while they performed a color
237 discrimination task at the center of the screen to ensure fixation (24). To identify and map
238 category-selective regions in the ventral temporal cortex, we presented subjects with four 5-
239 minute runs in which stimuli drawn from 5 categories (characters, bodies, faces, places, objects)
240 were presented in a block design (8 images per 4 second block), while subjects performed a 1-
241 back working memory task (25). Finally, to compare the neural response to natural images to
242 their synthesized scrambled counterparts, we presented subjects with at least eight 6-minute
243 runs, in which images were presented for 4 seconds, with no interstimulus interval, in an event-
244 related design (26, 27), while subjects performed the central fixation task described above.
245

246 Images subtended 12 degrees and were presented on both the left and right sides of the screen,
247 centered at an eccentricity of 7 degrees. We selected 10 different image classes, consisting of 7
248 objects and 3 textures, and for each class, presented 1 natural image and 2 synthesized images,
249 generated at a spatial constraint of 1x1 and from the pool4 layer. We also matched the Fourier
250 magnitude spectrum and the luminance histogram of the synthesized images to their
251 corresponding natural image, to control for potential low-level confounds. Over the course of the
252 entire experiment, each image was repeated approximately 20 to 24 times.

253 All scans were collected on a 3 Tesla General Electric MRI scanner, using a T2*
254 weighted sequence with multiplex factor of 4 (13 slices at multiplex 4 = 52 slices total), voxel size
255 of 2.5mm, repetition time (TR) of 1.0s and echo time (TE) of 30ms. Additionally, we acquired a
256 whole-brain high-resolution T1-weighted 3D BRAVO sequence with 0.9mm isotropic voxels. This
257 anatomical image was used for segmentation and surface reconstruction, which were performed
258 using Freesurfer. To correct for susceptibility distortions, we acquired an additional T2* weighted
259 sequence with reversed phase encoding direction and used the TOPUP function from FSL (28).
260 We performed volume-by-volume image registration to correct for motion artefacts using standard
261 procedures for motion correction (29). In the second session, we acquired another T1-weighted
262 3D BRAVO scan with voxel size 1.2 x 1.2 x 0.9mm. Using an image-based registration algorithm
263 (29), we aligned this anatomical scan to the high-resolution anatomical scan so that functional
264 regions of interest defined from the first session could be used to analyze the second session's
265 functional data.

266

267 **Defining cortical areas**

268 Using a 3-parameter population receptive field (pRF) model, we estimated the center
269 (x,y) and width (sigma) of the receptive field of each voxel in the occipital lobe (22). Then, we
270 manually drew visual area boundaries delineated by the reversal in the gradient of the polar angle
271 of pRFs (30). We were able to identify V1, V2, V3, and hV4 in all 7 subjects.

272 To identify category selective visual areas, we used the fLoc functional localizer (25), in
273 which images of faces, bodies, places, characters, objects, and phase-scrambles were presented
274 in a block design. We then used a GLM to estimate the response amplitudes to each stimulus
275 category and then performed a statistical contrast to identify category-selective voxels. We were
276 able to identify 3 face-selective clusters of voxels, in the mid-fusiform sulcus (mFus), posterior
277 fusiform gyrus (pFus) and inferior occipital gyrus (IOG) and 2 place-selective clusters of voxels, in
278 the transverse occipital sulcus (TOS) and the collateral sulcus (CoS), in each subject. We also
279 used an atlas-based approach to identify anatomically defined areas (31), using a surface-based
280 alignment to align the atlas to each subject's individual brain. We analyzed responses in 4 visual
281 areas from the Glasser Atlas: lateral occipital complex (LO), for which we combined 3 smaller
282 subregions, LO1, LO2, and LO3; ventral visual cortex (VVC); posterior inferotemporal cortex
283 (PIT); and ventromedial visual area (VMV) for which we combined VMV1, VMV2, and VMV3 (31).
284 These areas were selected because they have been identified as regions that contain information
285 about visual object category.

286 In all analyses, we thus examined a total of 13 visual areas: 4 retinotopically defined
287 areas (V1, V2, V3, hV4), 5 category-selective areas defined by a functional localizer (mFus, pFus,
288 IOG, TOS, CoS), and 4 anatomically defined areas from the Glasser atlas (LO, VVC, PIT, VMV).

289

290 **BOLD Data Analysis**

291 We extracted trial-averaged neural responses to individual images using the GLMdenoise
292 Matlab package (27), which estimates noise regressors from task-irrelevant voxels and uses
293 those in a generalized linear model (GLM) (32). To identify the most reliable voxels in each
294 cortical area, we split the data into two sets, such that each set contained half of the trials in
295 which a particular image was presented. Then, we re-fit the GLM separately to each of the two
296 splits and for each voxel, computed the correlation between its responses across the two splits.
297 This measure of split-half correlation was used to identify the most reliable voxels, and in all
298 analyses, we selected the 100 most reliable voxels in each ROI.

299 To determine how selective each visually responsive region is for the particular spatial
300 arrangement of features that is found in the natural image, we computed the Pearson distance
301 between the cortical response to a natural image and the neural response to a synth of the same

class ($d_{natural,synth}$). We also computed the Pearson distance between the cortical response to two different synths of the same class ($d_{synth,synth}$) (Fig. 4A). Finally, we computed the average Pearson distance between the cortical response to a synth of one class and the synths of every other class and called this the “between-class” distance. We assessed the category selectivity of a given cortical population by the degree to which the between-class distance exceeded the within class distance.

Triangle plot visualization

To visualize the relative representational distances between pairs of images, we plotted images in a triangle, where the length of the edges represents the magnitude of the representational distance between that pair of images. The representational distances are computed as the Pearson distance between each pair of images, for the dCNNs, cortical responses, and the model IT responses, but are estimated using maximum likelihood estimation for the human perceptual distances. Given the distances between 3 images, it is always possible to create a triangle where the edges correspond to distances, as long as none of the edge lengths exceeds the sum of the other two edge lengths. Then we rotate and translate the triangle so that the oddity image is always placed at the origin and the non-oddy images are above the oddity.

For the human BOLD responses, we also estimated the split-half distance as a measure of reliability of the response. To do so, we split all the trials for each subject into two halves and separately estimated the responses using the GLM for each half, then estimated the Pearson distance between the multivariate response to an image for one half compared to another half of the data. This split half distance is visualized in the triangle plots as a gray cloud around each image (Fig. 6D, 6I).

Readout analyses

Selectivity index

We quantified the selectivity for natural feature arrangement by the degree to which the natural-synth distance exceeded the synth-synth distance, normalized by the sum of the natural-synth distance and the synth-synth distance.

$$Selectivity\ Index = \frac{d_{natural,synth} - d_{synth1,synth2}}{d_{natural,synth} + d_{synth1,synth2}}$$

This measure reflects the extent to which a representation differentiates the natural image, i.e. the extent to which the natural image is more different from the synthesized images than the synthesized images are from each other.

Image-general readout

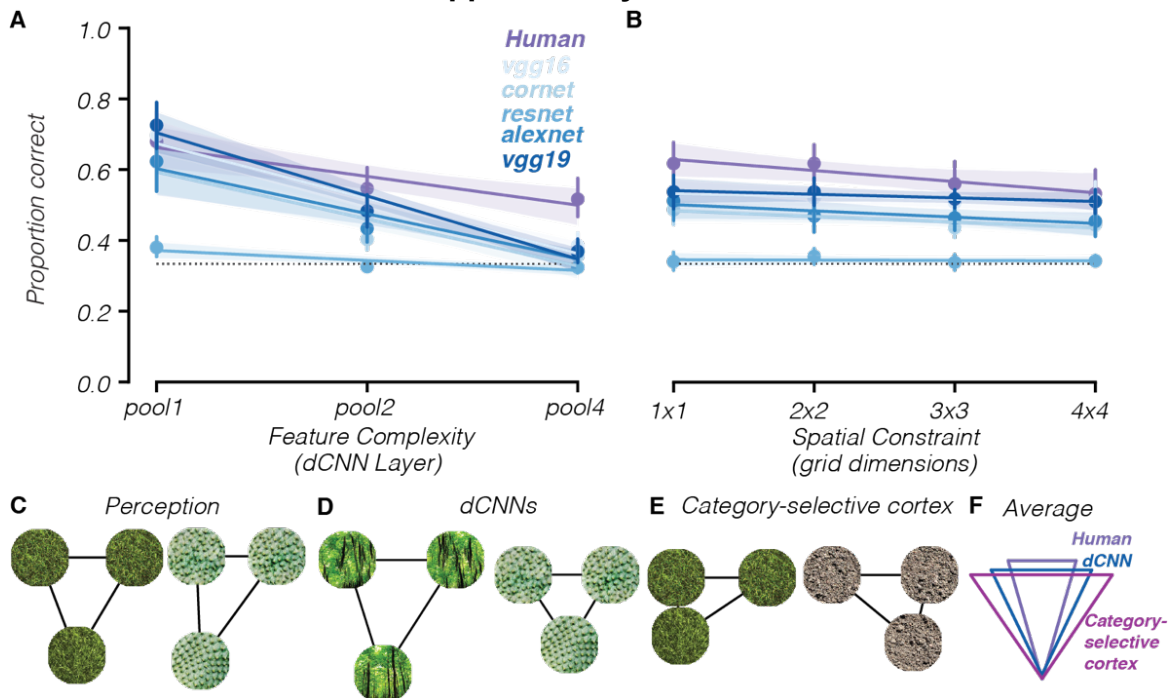
We performed an image-general readout by fitting weights to each voxel with the objective of maximizing the selectivity index across all image classes. We tested for generalization by fitting the weights to maximize the selectivity index for all but one image class and then evaluating the selectivity index on the held-out image class. Therefore, the number of parameters was equal to 100 (number of voxels) per area.

Image-specific readout

We performed an image class-specific readout by fitting a separate set of weights to each voxel for each image class, with the objective of maximizing the selectivity index for each image class separately. Therefore, this approach required 1000 parameters per visual area (10 voxels x 10 images). To prevent overfitting, we estimated betas for each trial separately, then randomly selected 90% of the trials, averaged together the betas, and fit the weights on that 90% of trials for each image class separately. Then we evaluated the selectivity index on the held-out 10% of trials. We selected 100 voxels for inclusion in this analysis by separately splitting up the 90% of trials into two halves and choosing the voxels which had the highest split-half reliability in this subset of the data. This approach therefore ensured that no part of the weight estimation could be influenced by the held-out trials.

356

Supplementary Results



357

358

359

360

361

362

363

364

365

366

367

368

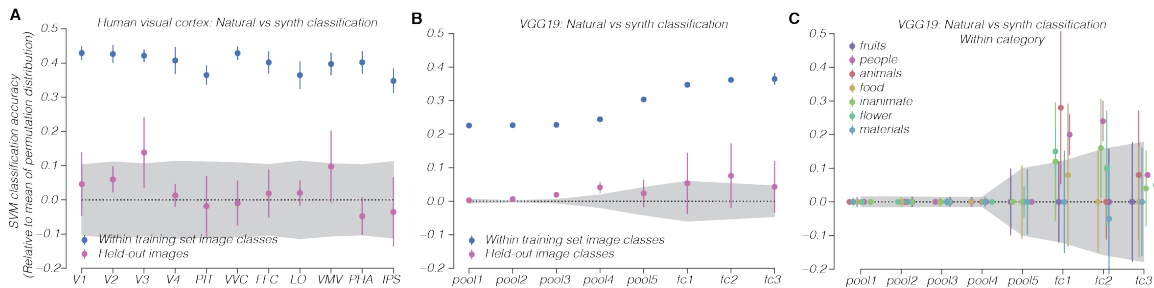
Supp. Fig. S1. Human observers are less sensitive to natural feature arrangement for texture-like images, similar to dCNN observer models and VTC voxels. (A) Performance of human observers (purple) compared to dCNN observer models (blues) at identifying the natural image as a function of feature complexity of synthesized images. (B) Performance of human observers (purple) compared to dCNN observer models (blues) at identifying the natural image as a function of constraints on spatial arrangements. (C-E) Triangular distance plots for perception (C), dCNN observer models (D), and category-selective cortical areas (E). (F) Average triangular distance plot across image categories, comparing category-selective cortex (magenta), dCNNs (blue), human observers (purple).

369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393

Classification Analyses

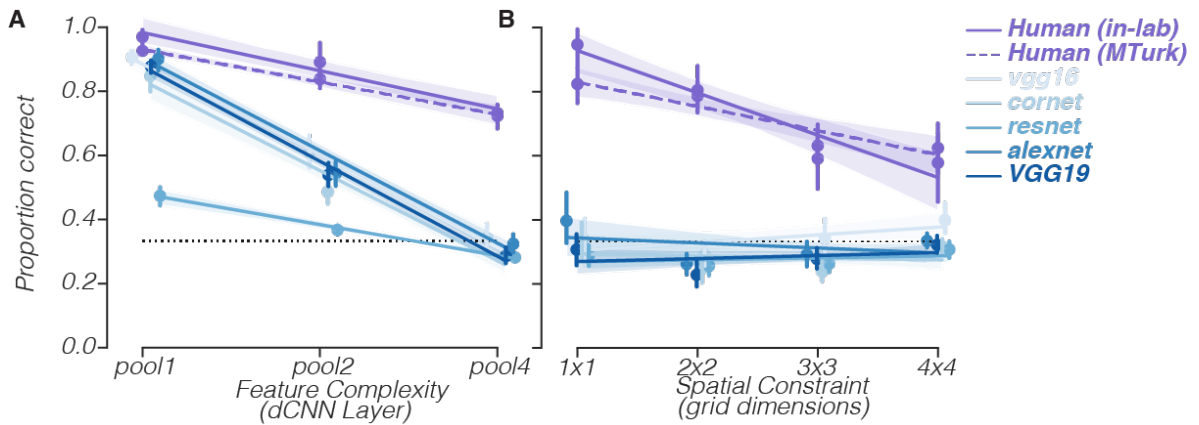
Using a support vector machine (SVM) classifier trained to classify images as natural or synthesized, we found corroborating evidence that cortical responses contain sufficient information to distinguish natural from synthesized images, though not in a generalizable format. We trained a SVM classifier with a linear kernel on cortical responses from 13 different visual areas. When evaluated on samples from image classes within its training set, the classifier was highly accurate in classifying the sample as natural or synthesized (Fig. S2A, blue points). However, when evaluated on samples from image classes outside its training set, the classifier was unable to classify the images as natural or synthesized significantly above chance level, computed using a permutation test (Fig. S2A, magenta points). Thus, a linear classification boundary can be found that distinguishes natural from scrambled images, but the classification boundary varies for different image classes.

To address the possibility that information about natural feature arrangement is present in the dCNN representation but dominated by information about unlocalized features, we trained a support vector machine (SVM) classifier with a linear kernel to predict whether an image was natural or synthesized. We found that when the SVM classifier was evaluated on image classes that were not within its training set, it was unable to predict whether these images were natural or synthesized significantly above chance (Fig. S2B, magenta points), even if the SVM was trained exclusively on other image classes within the same category (Fig. S2C). However, when evaluated on image classes within its training set (Fig. S2B, blue points), the SVM classifier was able to predict whether an image was natural or synthesized. These results suggest that the representation of natural and synthesized images is sufficiently different that an image-specific classification boundary can be found but not an image-general boundary.



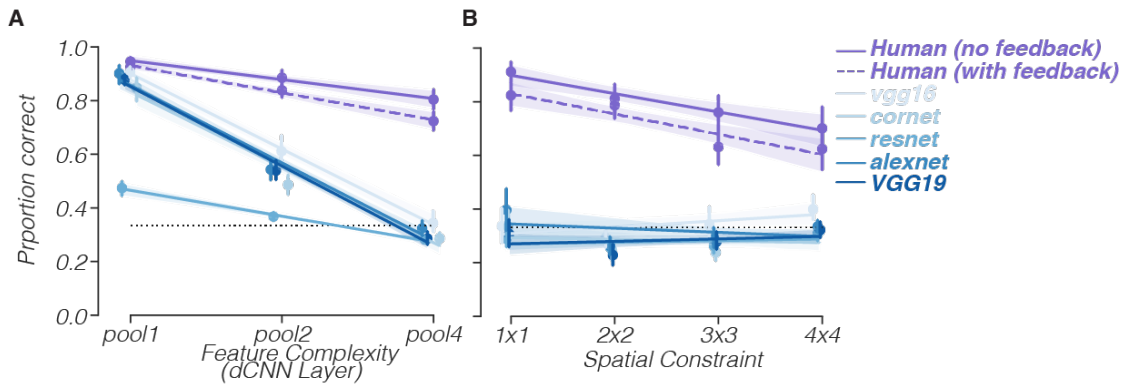
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408

Supp. Fig. S2. Classification accuracy using support vector machines to classify natural vs synth. (A) Human visual cortex natural-vs-synth classification accuracy, relative to mean of permutation distribution. Gray shaded region represents 95% confidence interval of permutation distribution. Blue points are classification accuracy for image classes within the training set, and pink points are classification accuracy for image classes that were not in the training set. (B) Same as A but using features from various VGG19 layers instead of cortical responses. (C) Within-category decoding accuracy. We grouped 37 image classes into 7 categories and trained a SVM classifier to predict whether an image was natural or synthesized on all image classes of the same category except one and evaluated its performance on the held out image class of the same category. Across 7 categories (fruits, people, animals, food, flowers, inanimate objects, and materials), we found that classification accuracy failed to exceed chance in layers pool1, pool2, pool3, pool4, and pool5, although classification accuracy did exceed chance level for two categories (animals, people) in fc1 and one category (people) in fc2.



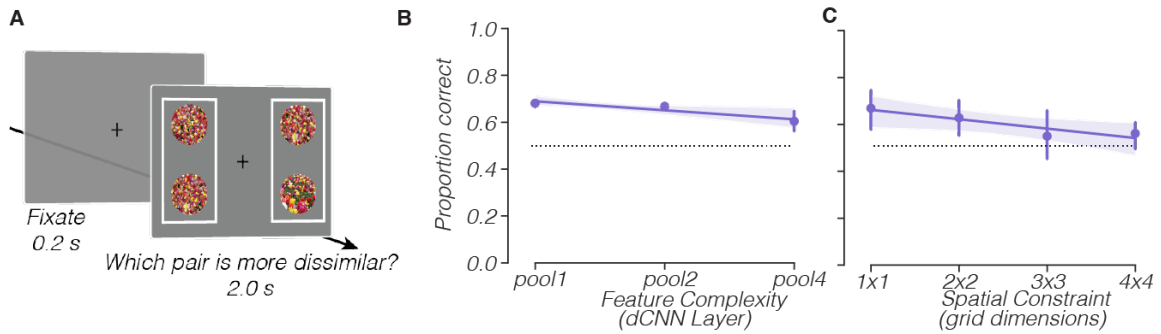
409
 410
 411
 412
 413
 414
 415
 416
 417

Supp. Fig. S3. Replication of behavioral results using dataset collected in-lab where fixation could be enforced with eye-tracking. (A) Comparison of human and dCNN behavior as a function of feature complexity. Solid purple line represents in-lab data and dashed purple line represents online data. (B) Comparison of human and dCNN behavior as a function of spatial constraint, fixing the feature complexity at the highest level (pool4).



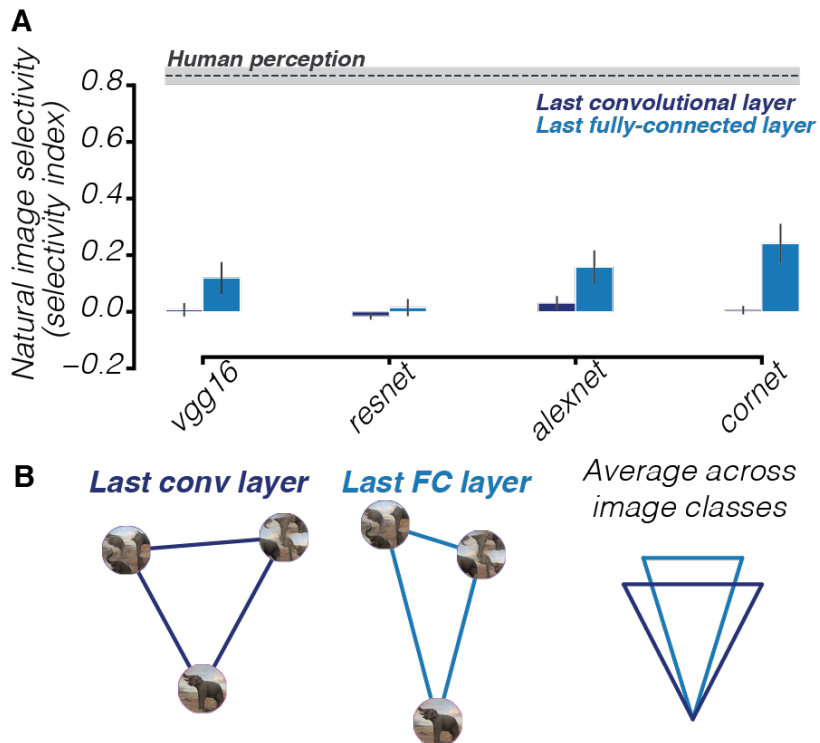
418
 419
 420
 421
 422
 423
 424
 425
 426
 427

Supp. Fig. S4. Replication of behavioral results using dataset collected on MTurk without correct/incorrect feedback. (A) Comparison of human behavior with feedback (solid purple line), human behavior without feedback (dashed purple line), and dCNN behavior (blue lines) as a function of feature complexity. (B) Comparison of human behavior with feedback, human behavior without feedback, and dCNN behavior as a function of spatial constraint, fixing the feature complexity at the highest level (pool4).



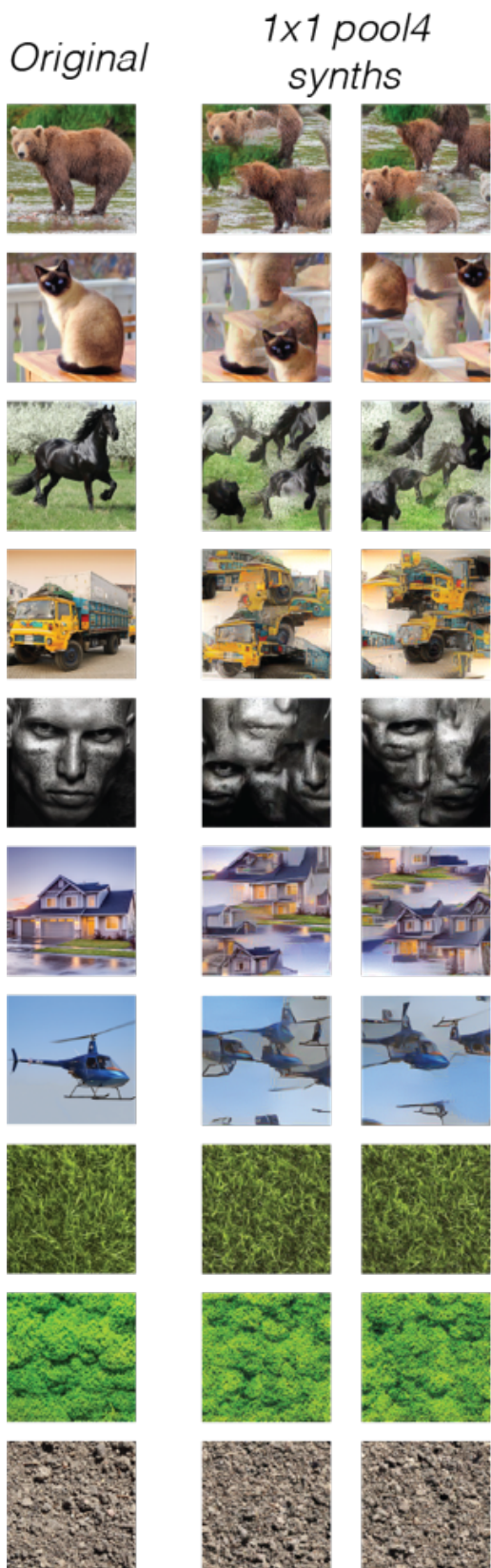
428
429
430
431
432
433
434
435
436
437

Supp. Fig. S5. Human behavior as a function of feature complexity and spatial constraint for the pairwise dissimilarity judgment task. (A) Task design. Subjects were shown two pairs of images and asked to select the pair which was more dissimilar from each other. (B) Human behavior as a function of feature complexity. The proportion of trials where subjects chose the pair with the natural image declined as the synths had more complex visual features. (C) Human behavior as a function of spatial constraint. The proportion of trials where subjects chose the pair with the natural declined as the arrangement of features in the synths was more strongly spatially constrained.



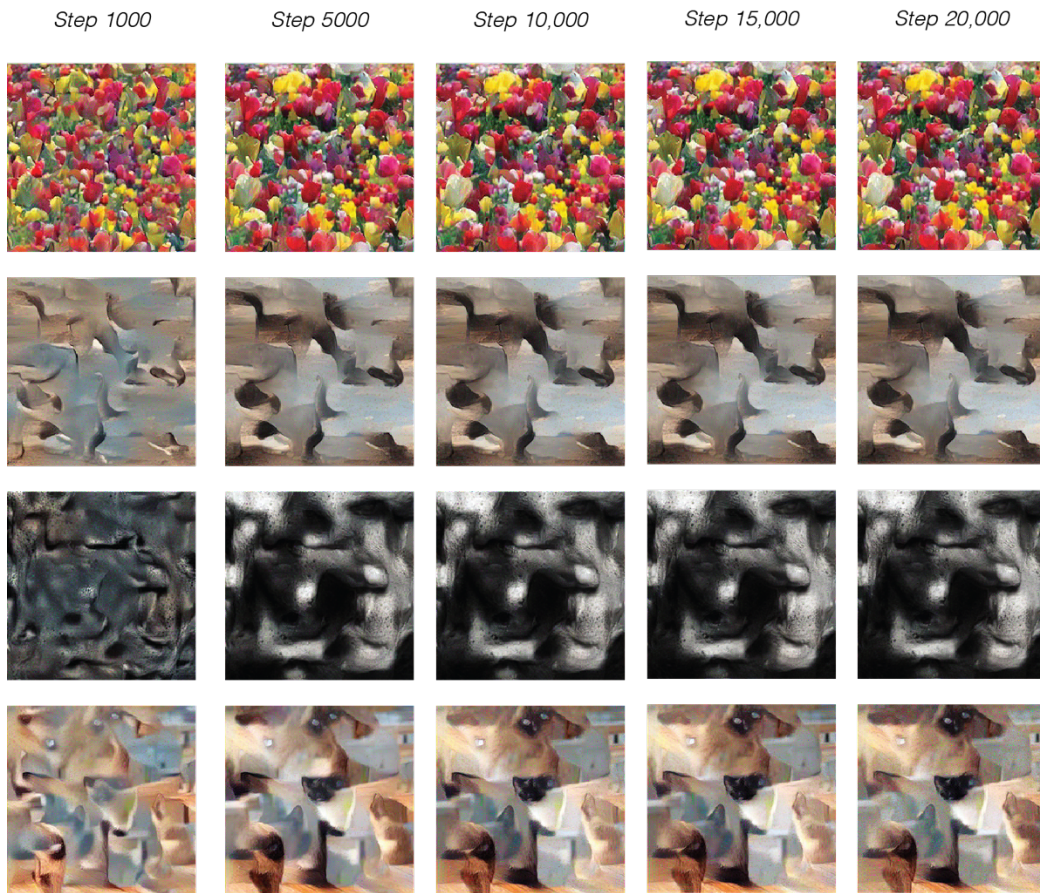
438
439
440
441
442
443

Supp. Fig. S6. Natural image selectivity for different dCNNs, comparing last convolutional layer to last fully-connected layer. (A) In all but one dCNN, selectivity for natural feature arrangement increases from the last convolutional layer to the last fully-connected layer. (B) Representational geometry comparing last convolutional layer to last fully connected layer.



444
445
446

Supp. Fig. S7. All stimuli used in neuroimaging experiment: 10 image classes consisting of 1 natural image and 2 synths (1x1 pool4 condition) per image class.



447
448
449
450
451

Supp. Fig. S8. Examples of synthesized images at different numbers of iterations in the synthesis process.

452 **Extended References**

453

- 454 1. L. Gatys, A. S. Ecker, M. Bethge, Texture synthesis using convolutional neural networks.
455 *Advances in neural information processing systems* **28**, 262–270 (2015).
- 456 2. K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image
457 Recognition. *Arxiv* (2014).
- 458 3. D. C. Liu, J. Nocedal, On the limited memory BFGS method for large scale optimization. *Math*
459 *Program* **45**, 503–528 (1989).
- 460 4. J. Portilla, E. P. Simoncelli, A Parametric Texture Model Based on Joint Statistics of Complex
461 Wavelet Coefficients. *Int J Comput Vision* **40**, 49–70 (2000).
- 462 5. J. Freeman, E. P. Simoncelli, Metamers of the ventral stream. *Nat Neurosci* **14**, 1195–201
463 (2011).
- 464 6. U. Güçlü, M. A. J. van Gerven, Deep Neural Networks Reveal a Gradient in the Complexity of
465 Neural Representations across the Ventral Stream. *J Neurosci* **35**, 10005–10014 (2015).
- 466 7. C. Olah, A. Mordvintsev, L. Schubert, Feature Visualization. *Distill* **2** (2017).
- 467 8. S. A. Cadena, *et al.*, Deep convolutional models improve predictions of macaque V1 responses
468 to natural images. *Plos Comput Biol* **15**, e1006897 (2019).
- 469 9. M. D. Zeiler, R. Fergus, Computer Vision – ECCV 2014, 13th European Conference, Zurich,
470 Switzerland, September 6-12, 2014, Proceedings, Part I. *Lect Notes Comput Sc*, 818–833 (2014).
- 471 10. N. C. L. Kong, B. Kaneshiro, D. L. K. Yamins, A. M. Norcia, Time-resolved correspondences
472 between deep neural network layers and EEG measurements in object processing. *Vision Res*
473 **172**, 27–45 (2020).
- 474 11. E. P. Simoncelli, W. T. Freeman, The steerable pyramid: a flexible architecture for multi-scale
475 derivative computation. *Proc Int Conf Image Process* **3**, 444–447 vol.3 (1995).
- 476 12. D. P. Kingma, J. Ba, Adam: A Method for Stochastic Optimization. *Arxiv* (2014).
- 477 13. M. Abadi, *et al.*, TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed
478 Systems. *Arxiv* (2016).
- 479 14. T. S. A. Wallis, *et al.*, A parametric texture model based on deep convolutional features
480 closely matches texture appearance for humans. *J Vision* **17**, 5 (2017).
- 481 15. P. Virtanen, *et al.*, SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat*
482 *Methods* **17**, 261–272 (2020).
- 483 16. N. J. Majaj, H. Hong, E. A. Solomon, J. J. DiCarlo, Simple Learned Weighted Sums of Inferior
484 Temporal Neuronal Firing Rates Accurately Predict Human Core Object Recognition
485 Performance. *J Neurosci Official J Soc Neurosci* **35**, 13402–18 (2015).
- 486 17. D. L. K. Yamins, *et al.*, Performance-optimized hierarchical models predict neural responses
487 in higher visual cortex. *P Natl Acad Sci Usa* **111**, 8619–24 (2014).
- 488 18. M. Schrimpf, *et al.*, Brain-Score: Which Artificial Neural Network for Object Recognition is
489 most Brain-Like? *Biorxiv*, 407007 (2020).
- 490 19. D. L. K. Yamins, J. J. DiCarlo, Using goal-driven deep learning models to understand sensory
491 cortex. *Nat Neurosci* **19**, 356–65 (2016).
- 492 20. B. A. Richards, *et al.*, A deep learning framework for neuroscience. *Nat Neurosci* **22**, 1761–
493 1770 (2019).
- 494 21. S. Ogawa, T. M. Lee, A. R. Kay, D. W. Tank, Brain magnetic resonance imaging with contrast
495 dependent on blood oxygenation. *Proc National Acad Sci* **87**, 9868–9872 (1990).
- 496 22. S. O. Dumoulin, B. A. Wandell, Population receptive field estimates in human visual cortex.
497 *Neuroimage* **39**, 647–660 (2008).
- 498 23. B. A. Wandell, J. Winawer, Imaging retinotopic maps in the human brain. *Vision Res* **51**, 718–
499 737 (2011).
- 500 24. J. L. Gardner, E. P. Merriam, J. A. Movshon, D. J. Heeger, Maps of visual space in human
501 occipital cortex are retinotopic, not spatiotopic. *J Neurosci Official J Soc Neurosci* **28**, 3988–99
502 (2008).
- 503 25. A. Stigliani, K. S. Weiner, K. Grill-Spector, Temporal Processing Capacity in High-Level Visual
504 Cortex Is Domain Specific. *J Neurosci* **35**, 12412–12424 (2015).
- 505 26. B. Vintch, J. L. Gardner, Cortical correlates of human motion perception biases. *J Neurosci*
506 *Official J Soc Neurosci* **34**, 2592–604 (2014).

507 27. K. N. Kay, A. Rokem, J. Winawer, R. F. Dougherty, B. A. Wandell, GLMdenoise: a fast,
508 automated technique for denoising task-based fMRI data. *Front Neurosci-switz* **7**, 247 (2013).
509 28. J. L. R. Andersson, S. Skare, J. Ashburner, How to correct susceptibility distortions in spin-
510 echo echo-planar images: application to diffusion tensor imaging. *Neuroimage* **20**, 870–888
511 (2003).
512 29. O. Nestares, D. J. Heeger, Robust multiresolution alignment of MRI brain volumes. *Magnet*
513 *Reson Med* **43**, 705–715 (2000).
514 30. J. Larsson, D. J. Heeger, Two Retinotopic Visual Areas in Human Lateral Occipital Cortex. *J*
515 *Neurosci* **26**, 13128–13142 (2006).
516 31. M. F. Glasser, *et al.*, A multi-modal parcellation of human cerebral cortex. *Nature* **536**, 171–
517 178 (2016).
518 32. K. J. Friston, *et al.*, Statistical parametric maps in functional imaging: A general linear
519 approach. *Hum Brain Mapp* **2**, 189–210 (1994).
520 33. D. Birman, J. L. Gardner, A quantitative framework for motion visibility in human cortex. *J*
521 *Neurophysiol* **120**:1824–1839 (2018).