# JNeurosci
## THE JOURNAL OF NEUROSCIENCE

*Research Articles: Behavioral/Cognitive*

# Categorical Biases in Human Occipitoparietal Cortex

# Categorical Biases in Human Occipitoparietal Cortex

Edward F. Ester[1*], Thomas C. Sprague[2], John T. Serences[3]

[1]*Department of Psychology, Center for Complex Systems and Brain Sciences, and FAU Brain Institute, Florida Atlantic University*
[2]*Department of Psychological and Brain Sciences, University of California, Santa Barbara*
[3]*Department of Psychology, Neurosciences Graduate Program, and Kavli Institute for Brain and Mind, University of California, San Diego*

*Correspondence:

Edward Ester
Department of Psychology, Center for Complex Systems and Brain Sciences, and FAU Brain Institute
Florida Atlantic University
777 Glades Rd.
Boca Raton, FL. 33431
eester@fau.edu

**Word Counts:**

Abstract: 183/250
Significance Statement: 93/120
Introduction: 463/650
Discussion: 1190/1500

Figure/Table Count: 18/0

**Abstract**

35

36        Categorization allows organisms to generalize existing knowledge to novel stimuli and to

37    discriminate between physically similar yet conceptually different stimuli. Humans, nonhuman

38    primates, and rodents can readily learn arbitrary categories defined by low-level visual features,

39    and learning distorts perceptual sensitivity for category-defining features such that differences

40    between physically similar yet categorically distinct exemplars are enhanced while differences

41    between equally similar but categorically identical stimuli are reduced. We report a possible

42    basis for these distortions in human occipitoparietal cortex. In three experiments, we used an

43    inverted encoding model to recover population-level representations of stimuli from multivoxel

44    and multi-electrode patterns of human brain activity while human participants (both sexes)

45    classified continuous stimulus sets into discrete groups. In each experiment, reconstructed

46    representations of to-be-categorized stimuli were systematically biased towards the center of the

47    appropriate category. These biases were largest for exemplars near a category boundary,

48    predicted participants' overt category judgments, emerged shortly after stimulus onset, and could

49    not be explained by mechanisms of response selection or motor preparation. Collectively, our

50    findings suggest that category learning can influence processing at the earliest stages of cortical

51    visual processing.

52

53

**Significance Statement**

Category learning enhances perceptual sensitivity for physically similar yet categorically different stimuli. We report a possible mechanism for these changes in human occipitoparietal cortex. In three experiments, we used an inverted encoding model to recover population-level representations of stimuli from multivariate patterns in occipitoparietal cortex while participants categorized sets of continuous stimuli into discrete groups. The recovered representations were systematically biased by category membership, with larger biases for exemplars adjacent to a category boundary. These results suggest that mechanisms of categorization shape information processing at the earliest stages of the visual system.

65      Categorization refers to the process of mapping continuous sensory inputs onto discrete

66  and behaviorally relevant concepts. It is a cornerstone of flexible behavior that allows organisms

67  to generalize existing knowledge to novel stimuli and to discriminate between physically similar

68  yet conceptually different stimuli. Many real-world categories are defined by a combination of

69  low-level visual properties such as hue, luminance, spatial frequency, and orientation. For

70  example, a forager might be tasked with determining whether a food source is edible based on

71  subtle variations in color, shape, size, and texture. Humans and other animals can readily learn

72  arbitrary novel categories defined by low-level visual properties (Goldstone, 1998; Ashby &

73  Maddox, 2005), and such learning "distorts" perceptual sensitivity for category-defining features

74  such that discrimination performance for physically similar yet categorically different stimuli is

75  increased (i.e., acquired distinctiveness; Goldstone, 1995; Newell & Bulthoff, 2002) and

76  discrimination performance for stimuli from the same category reduced (i.e., acquired similarity;

77  Livingston et al., 1998).

78      Invasive electrophysiological studies suggest that single-unit responses in early visual

79  areas index the physical properties of a stimulus but not its category membership, while single-

80  unit responses in later areas index the category membership of a stimulus regardless of its

81  physical properties (e.g., Sigala & Logothetis, 2002; Freedman et al., 2001; Freedman & Assad,

82  2006). These results have been taken as evidence that category-selective responses are a *de novo*

83  property of higher-order visual areas. However, perceptual distortions following category

84  learning could also reflect subtle changes in how to-be-categorized information is represented by

85  sensory neural populations (Folstein et al., 2012; Davis & Poldrack, 2013). Here we provide a

86  test of this possibility. In three experiments, we trained human participants (both sexes) to

87  classify sets of continuous stimuli into discrete groups. Next, next, we applied multivariate

88    models to noninvasive measurements of human brain activity (fMRI and EEG) from visual and

89    parietal cortical areas while participants categorized the same stimulus sets. This allowed us to

90    recover, visualize, and quantify stimulus-specific representations of to-be-categorized exemplars.

91    In Experiment 1 (fMRI), we show that reconstructed representations of to-be-categorized

92    orientations in visual areas V1-V3 are systematically biased towards the center of the category to

93    which they belong. These biases were correlated with trial-by-trial variability in overt category

94    judgments and were largest for orientations adjacent to the category boundary where they would

95    be most beneficial for discrimination performance. In Experiment 2, we utilized EEG to generate

96    time-resolved representations of to-be-categorized orientations and show that categorical biases

97    manifest shortly after stimulus onset ($\leq 300$ ms). In Experiment 3, we used EEG and a delayed

98    match-to-category task to show that categorical biases observed in Experiments 1 and 2 cannot

99    be explained by response biases or motor preparation. Collectively, our findings suggest that

100   mechanisms of categorization can shape information processing at the earliest stages of the

101   visual system.

102

**Methods**

103

*General Overview*

104

*Participants*. A total of 44 human volunteers (both sexes) participated in this study. Eight

105

participants completed Experiment 1 (fMRI), 28 participants completed Experiment 2 (EEG),

106

and eight participants completed Experiment 3 (EEG). Experiments 1 and 2 were performed at

107

the University of California, San Diego, while Experiment 3 was performed at Florida Atlantic

108

University. Participants were recruited from the student body at each university. All study

109

procedures were approved by local institutional review boards, and all participants gave both

110

written and oral informed consent. Participants self-reported normal or corrected-to-normal

111

visual acuity and were remunerated with cash incentives ($20/h for fMRI and $15/h for EEG).

112

*Stimulus Displays*. Stimulus displays were generated in MATLAB and rendered using

113

Psychophysics Toolbox software extensions (Kleiner et al., 2017). During Experiment 1 (fMRI),

114

displays were projected onto a 110 cm-wide screen placed at the base of the MRI table, and

115

participants viewed displays via a mirror attached to the MR head coil from a distance of 370

116

cm. During Experiments 2 and 3, displays were projected onto a 19-inch CRT monitor cycling at

117

120Hz (Experiment 2) or 85Hz (Experiment 3). Participants were seated approximately 65 cm

118

from the display (head position was not constrained).

119

*Experiment 1 - fMRI*

120

*Participants*. Eight neurologically intact human volunteers (AA, AB, AC, AD, AE, AF, AG, and

121

AH; six females) completed Experiment 1. Each participant completed a single one-hour

122

behavioral training session approximately 24-72 hours prior to scanning. Seven participants (AA,

123

AB, AC, AD, AE, AF, AG) completed two 2-hour experimental scan sessions; an eighth

124

participant (AH) completed a single 2-hour experimental scan session. Participants AA, AB, AC,

125

6

126    AD, AE, AF, and AH also completed a single 2-hour retinotopic mapping scan session. Data

127    from this session were used to identify visual field borders in early visual cortical areas V1-

128    hV4/V3A and subregions of posterior intraparietal sulcus (IPS0-3; see *Retinotopic Mapping*,

129    below).

130    *Behavioral Tasks*. In separate runs (where "run" refers to a continuous block of 30 trials lasting

131    280 seconds) participants performed either an orientation mapping task or a category

132    discrimination task. Trials in each task lasted 3 seconds, and consecutive trials were separated by

133    a 5 or 7 s inter-trial-interval (pseudorandomly chosen on each trial). During the orientation

134    mapping task, participants attended a stream of letters presented at fixation (subtending 1.0˚ x

135    1.0˚ from a viewing distance of 370 cm) while ignoring a task-irrelevant phase-reversing (15 Hz)

136    square-wave grating (0.8 cycles/deg with inner and outer radii of 1.16˚ and 4.58˚, respectively)

137    presented in the periphery. On each trial, the grating was assigned one of 15 possible orientations

138    (0˚-168˚ in 12˚ increments). Participants were instructed to detect and report the identity of a

139    target ("X" or "Y") in the letter stream using an MR-compatible button box. Only one target was

140    presented on each trial. Letters were presented at a rate of 10 Hz (50% duty cycle, i.e. 50 msec

141    on, 50 msec off), and targets could occur during any cycle from +750 to +2250 msec after

142    stimulus onset. During category discrimination runs, participants were shown displays containing

143    a circular aperture (inner and outer radii of 1.16˚ and 4.58˚ from a viewing distance of 370 cm)

144    filled with 150 iso-oriented bars (see Figure 1A). Each bar subtended 0.2˚ x 0.6˚ with a stroke

145    width of 8 pixels (1024 x 768 display resolution). Each bar flickered at 30 Hz and was randomly

146    replotted within the aperture at the beginning of each "up" cycle.

147        On each trial, all bars were assigned an orientation from 0˚-168˚ in 12˚ increments.

148    Inspired by earlier work in non-human primates (Freedman & Assad, 2006), we randomly

7

149     selected and designated one of these orientations as a category boundary such that the seven

150     orientations counterclockwise to this value were assigned membership in "Category 1", while the

151     seven orientations clockwise to this value were assigned membership in "Category 2".

152     Participants were not informed that the category boundary was chosen from the set of possible

153     stimulus orientations. Participants reported whether the orientation shown on each trial was a

154     member of Category 1 or 2 (via an MR-compatible button box). Participants were free to respond

155     at any point during the trial, though the stimulus was always presented for a total of 3000 ms.

156     Each participant was familiarized and trained to criterion performance on the category

157     discrimination task during a one-hour behavioral testing session completed one to three days

158     prior to his or her first scan session. Written feedback ("Correct!" or "Incorrect") was presented

159     in the center of the display for 1.25 sec. after each trial during behavioral training and MR

160     scanning. Across either one (N = 1) or two (N = 7) scan sessions, each participant completed 7

161     (N = 1), 13 (N = 1), 14 (N = 1), 15 (N = 1) or 16 (N = 4) runs of the orientation mapping and

162     category discrimination tasks.

163     *fMRI Acquisition and Preprocessing*. Imaging data were acquired with a 3.0T GE MR 750

164     scanner located at the Center for Functional Magnetic Resonance imaging on the UCSD campus.

165     All images were acquired with a 32 channel Nova Medical head coil (Wilmington, MA). Whole-

166     brain echo-planar images (EPIs) were acquired in 35 3 mm slices (no gap) with an in-plane

167     resolution of 3 x 3 mm (192 x 192 mm field-of-view, 64 x 64 mm image matrix, 90° flip angle,

168     2000 ms TR, 30 ms TE). During retinotopic mapping scans (see below) EPIs were acquired in 31

169     3mm thick oblique slices (no gap) positioned over posterior visual and parietal cortex with an in-

170     plane resolution of 2 x 2 mm (192 x 192 mm field-of-view, 96 x 96 mm image matrix, 90° flip

171     angle, 2250 ms TR, 30 ms TE). EPIs were coregistered to a high-resolution anatomical image

172    collected during the same session (FSPGR T1-weighted sequence, 11 ms TR, 3.3 ms TE, 1100

173    ms TI, 172 slices, 18° flip angle, 1 mm$^3$ resolution), unwarped (FSL software extensions), slice-

174    time-corrected, motion-corrected, high-pass-filtered (to remove first-, second- and third-order

175    drift), transformed to Talairach space, and normalized (z-score) on a scan-by-scan basis. Data

176    from data from scan sessions were then co-registered to a high-resolution anatomical image

177    (FSPGR T1-weighted sequences; parameters as described above) collected during the retinotopic

178    mapping session.

179    *Retinotopic Mapping*. Retinotopically organized visual areas V1-hV4v/V3A were defined using

180    data from a single retinotopic mapping run collected during each experimental scan session.

181    Participants fixated a small dot at fixation while phase-reversing (8 Hz) checkerboard wedges

182    subtending 60° of polar angle (at maximum eccentricity) were presented along the horizontal or

183    vertical meridian (alternating with a period of 40 seconds; i.e., 20 seconds of horizontal

184    stimulation followed by 20 seconds of vertical stimulation). To identify visual field borders, we

185    constructed a general linear model with two boxcar regressors, one marking epochs of vertical

186    stimulation and another marking epochs of horizontal stimulation. Each regressor was convolved

187    with a canonical hemodynamic function ("double gamma" as implemented in BrainVoyager

188    QX). Next, we generated a statistical parametric map marking voxels with larger responses

189    during epochs of vertical relative to horizontal stimulation. This map was projected onto a

190    computationally inflated representation of each participant's cortical surface for visualization to

191    aid in the definition of the borders of visual areas V1, V2v, V2d, V3v, V3d, hV4v, and V3A.

192    Data from V2v and V2d were combined into a single V2 ROI, and data from V3v and V3d were

193    combined into a single V3 ROI. ROIs were also combined across cortical hemispheres (e.g., left

194  and right V1) as no asymmetries were observed and the stimulus was presented in the center of

195  the visual field.

196       Seven participants (AA, AB, AC, AD, AE, AF, and AH) completed a separate two-hour

197  retinotopic mapping scan; data from this session were used to identify retinotopically organized

198  regions of inferior parietal sulcus (IPS0-3). During each task run, participants were shown

199  displays containing a rotating wedge stimulus (period 24.75 or 36 sec) that subtended 72˚ of

200  polar angle with inner and outer radii of 1.75 and 8.75˚, respectively. In alternating blocks, the

201  wedge contained a 4 Hz phase-reversing checkerboard or field of moving dots and participants

202  were required to detect small, brief, and temporally unpredictable changes in checkerboard

203  contrast or dot speed. Six participants completed between 8 and 14 task runs. To compute the

204  best polar angle for each voxel in IPS we shifted the signals from counterclockwise runs by twice

205  the estimated hemodynamic response function (HRF) delay (2 x 6.75 s = 13.5 s), removed data

206  from the first and last full stimulus cycle, and reversed the time series so that all runs reflected

207  clockwise rotation. We next computed the power and phase of the response at the stimulus'

208  period (either 1/24.75 or 1/36 Hz) and subtracted the estimated hemodynamic response function

209  delay (6.75 seconds) to align the signal phase in each voxel with the stimulus' location. Maps of

210  orientation preference (computed via cross-correlation) were projected onto a computationally

211  inflated representation of each participant's grey-white matter boundary to aide in the

212  identification of visual field borders separating IPS0-3. An eighth participant (AG) chose not to

213  participate in an additional retinotopic mapping session. For this participant, we estimated visual

214  field borders for visual areas V1-hV4/V3A. using data from the retinotopic mapping run

215  collected during the participant's sole experimental session. We did not attempt to define IPS

216  regions IPS0-3 for this participant.

217    *Decoding Categorical Biases in Visual Cortex*. We used a linear decoder to examine whether

218    fMRI activation patterns evoked by exemplars adjacent to the category boundary and at the

219    center of each category were more similar during the category discrimination task relative to the

220    orientation mapping task (i.e., acquired similarity). In the first phase of the analysis, we trained a

221    linear support vector machine (LIBSVM implementation; Chang & Lin, 2011) to discriminate

222    between the oriented exemplars at the center of each category (48° from the boundary) using

223    data from the orientation mapping and category discrimination tasks. To ensure internal

224    reliability, we implemented a "leave-one-run-out" cross validation scheme where data from all

225    but one scanning run was used to train the classifier and data from the remaining scanning run

226    were used for validation. This procedure was repeated until data from each scan had served as

227    the validation set, and the results were averaged across permutations. Next, we trained a second

228    classifier on activation patterns evoked by exemplars at the center of each category boundary and

229    used the trained classifier to predict the category membership of exemplars adjacent to the

230    category boundary. If category learning increases the similarity of activation patterns evoked by

231    exemplars within the same category, then within-category decoding performance should be

232    superior during the category discrimination task relative to the orientation mapping task.

233    *Inverted Encoding Model of Orientation Selectivity.* A linear inverted encoding model (IEM)

234    was used to recover a model-based representation of stimulus orientation from multivoxel

235    activation patterns measured in early visual areas (Brouwer & Heeger, 2011). The same general

236    approach was used during Experiments 1 (fMRI) and 2 (EEG). Specifically, we modeled the

237    responses of voxels (electrodes) measured during the orientation mapping task as a weighted

238    sum of 15 orientation-selective channels, each with an idealized response function (half-wave-

239    rectified sinusoid raised to the $14^{th}$ power). The maximum response of each channel was set to

240     unit amplitude; thus units of response are arbitrary. Let $B_1$ (*m* voxels or electrodes x $n_1$ trials) be

241     the response of each voxel (electrode) during each trial of the RSVP task, let $C_1$ (*k* filters x $n_1$

242     trials) be a matrix of hypothetical orientation filters, and let W (*m* voxels or electrodes x *k* filters)

243     be a weight matrix describing the mapping between $B_1$ and $C_1$:

244

$$B_1 = W C_1$$

245

246     In the first phase of the analysis, we computed the weight matrix W from the voxel-wise

247     (electrode-wise) responses in $B_1$ via ordinary least-squares:

248

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

249

250     Next, we defined a test data set $B_2$ (*m* voxels or electrodes x $n_2$ trials) using data from the

251     category discrimination task. Given W and $B_2$, a matrix of filter responses $C_2$ (*k* filters x *n* trials)

252     can be estimated via model inversion:

253

$$C_2 = (W^T W)^{-1} W^T B_2$$

254

255     $C_2$ contains the reconstructed response of each modeled orientation channel (the channel

256     response function; CRF) on each trial of the category discrimination task. This analysis can be

257     considered a form of model-based, directed dimensionality reduction where activity patterns are

258     transformed from their original measurement space (fMRI voxels; EEG electrodes) into a

259     modeled information space (orientation-selective channels). Importantly, results from this

260 method cannot be used to infer any changes in orientation tuning – or any properties of neural

261 responses - occurring at the single neuron level, and only assay the information content of large-

262 scale patterns of neural activity (Sprague et al., 2018) Additionally, while it is the case that

263 arbitrary linear transforms can be applied to the basis set, model weights, and reconstructed

264 channel response function (Gardiner & Liu, 2019), results are uniquely defined for a given model

265 specification (Sprague, Boynton & Serences, 2019). Trial-by-trial CRFs were multiplied by the

266 original basis set to recover a full 180-degree function, circularly shifted to a common center (0°)

267 and sorted by category membership so that any category bias would manifest as a clockwise shift

268 (i.e., towards the center of Category 2).

269 *Quantification of Bias in Orientation Representations*. To quantify categorical biases in

270 reconstructed model-based CRFs, these functions were fit with an exponentiated cosine function

271 of the form:

272

$$f(x) = \alpha\left(e^{k(\cos(\mu-x)-1)}\right) + \beta$$

273

274 where, $x$ is a vector of channel responses and $\alpha$, $\beta$, $k$ and $\mu$ correspond to the amplitude (i.e.,

275 signal over baseline), baseline, concentration (the inverse of bandwidth) and the center of the

276 function, respectively. Fitting was performed using a multidimensional nonlinear minimization

277 algorithm (Nelder-Mead).

278   Category biases in the estimated center of each construction ($\mu$) during the category

279 discrimination task were quantified via permutation tests. For a given visual area (e.g., V1) we

280 randomly selected (with replacement) stimulus reconstructions from eight of eight participants.

281 Specifically, we computed a "mean" reconstruction by randomly selecting (with replacement)

282 and averaging reconstructions from all participants. The mean reconstruction was fit with the

283 cosine function described above, yielding point estimates of α, β, *k*, and *μ*. This procedure was

284 repeated 1,000 times, yielding 1,000 element distributions of parameter estimates. We then

285 computed the proportion of permutations where a μ value less than 0 was obtained to obtain an

286 empirical *p*-value for categorical shifts in reconstructed representations.

287 *Searchlight Decoding of Category Membership.* We used a roving searchlight analysis (Ester et

288 al., 2015) to identify cortical regions beyond V1-V3 that contained category-specific

289 information. We defined a spherical neighborhood with a radius of 8 mm around each grey

290 matter voxel in the cortical sheet. We next extracted and averaged the normalized response of

291 each voxel in each neighborhood over a period from 4-8 seconds after stimulus onset (this

292 interval was chosen to account for typical hemodynamic lag of 4-6 seconds). A linear SVM

293 (LIBSVM implementation) was used to classify stimulus category using activation patterns

294 within each neighborhood. To classify category membership, we designated the three

295 orientations immediately counterclockwise to the category boundary (see Figure 1) as members

296 of Category 1 and the three orientations immediately clockwise of the boundary as members of

297 Category 2. We then trained our classifier to discriminate between categories using data from all

298 but one task run. The trained classifier was then used to predict category membership from

299 activation patterns measured during the held-out task run. This procedure was repeated until each

300 task run had been held out, and the results were averaged across permutations. Finally, we

301 repeated the same analysis using the three Category 1 and Category 2 orientations adjacent to the

302 second (orthogonal) category boundary (see Figure 1) and averaged the results across category

303 boundaries.

14

304    We identified neighborhoods encoding stimulus category using a leave-one-participant-

305    out cross validation approach (Esterman et al., 2010). Specifically, for each participant (e.g., AA)

306    we randomly selected (with replacement) and averaged classifier performance estimates from

307    each neighborhood from each of the remaining 7 volunteers (e.g., AB-AH). This procedure was

308    repeated 1000 times, yielding a set of 1000 classifier performance estimates for each

309    neighborhood. We generated a statistical parametric map (SPM) for the held-out participant that

310    indexed neighborhoods where classifier performance was greater than chance (50%) on 97.5% of

311    permutations (false-discovery-rate corrected for multiple comparisons across neighborhoods).

312    Finally, we projected each participant's SPM onto a computationally inflated representation of

313    his or her grey-white matter boundary and used Brain Voyager's "Create POIs from Map

314    Clusters" function with an area threshold of 25 mm$^2$ to identify ROIs supporting above-chance

315    category classification performance. Because of differences in cortical folding patterns, some

316    ROIs could not be unambiguously identified in all 8 participants. Therefore, across participants,

317    we retained all ROIs that were shared by at least 7 out of 8 participants. Finally, we extracted

318    multivoxel activation patterns from each ROI and computed model-based reconstructions of

319    channel response functions during the RSVP and category tasks using a leave-one-run-out cross-

320    validation approach. Specifically, we used data from all but one task run to estimate a set of

321    orientation weights for each voxel in each ROI. We then used these weights and activation

322    patterns measured during the held-out task run to estimate a channel response function, which

323    contains a representation of stimulus orientation. This procedure was repeated until each task run

324    had been held out, and the results were averaged across permutations. Note that each

325    participant's ROIs were defined using data from the remaining 7 participants. This ensured that

326 participant-level reconstructions were statistically independent of the searchlight method used to

327 define ROIs encoding category information.

328 *Within-participant Error Bars*. We report estimates of within-participant variability (e.g., ±1

329 S.E.M.) throughout the paper. These estimates discard subject variance (e.g., overall differences

330 in BOLD response amplitude) and instead reflect variance related to the subject by condition(s)

331 interaction term(s) (i.e., variability in estimated channel responses). We used the approach

332 described by Cousineau (2005): raw data (e.g., channel response estimates) were de-meaned on a

333 participant by participant basis, and the grand mean across participants was added to each

334 participant's zero-centered data. The grand mean-centered data were then used to compute

335 estimates of standard error.

336 **Experiment 2 - EEG**

337 *Participants*. 29 new volunteers recruited from the UC San Diego community completed

338 Experiment 2. All participants self-reported normal or corrected-to-normal visual acuity and

339 gave both written and oral informed consent as required by the local Institutional Review Board.

340 Each participant was tested in a single 2.5-3 hour experimental session (the exact duration varied

341 across participants depending on the amount of time needed to set up and calibrate the EEG

342 equipment). Unlike Experiment 1, participants were not trained on the categorization task prior

343 to testing. We adopted this approach in the hopes of tracking the gradual emergence of

344 categorical biases during learning. However, many participants learned the task relatively

345 quickly (within 40-60 trials), leaving too few trials to enable a direct analysis of this possibility.

346 Data from one participant were discarded due to a high number of EOG artifacts (over 35% of

347 trials); the data reported here reflect the remaining 28 participants.

348 *Behavioral Tasks*.

349     In separate runs (where "run" refers to a continuous block of 60 trials lasting approximately 6.5

350     minutes), participants performed orientation mapping and category discrimination tasks similar

351     to those used in Experiment 1. During both tasks a rapid series of letters (subtending 1.14° x

352     1.14° from a viewing distance of 55 cm) was presented at fixation, and an aperture of 150 iso-

353     oriented bars (subtending 0.5° x 1.2°) was presented in the periphery. The aperture of bars had

354     inner and outer radii of 1.96° and 9.13°, respectively. On each trial, the bars were assigned one of

355     15 possible orientations (again 0°-168° in 12° increments) and flickered at a rate of 30 Hz. Each

356     bar was randomly replotted within the aperture at the beginning of each "up" cycle. Letters in the

357     RSVP stream were presented at a rate of 6.67 Hz

358         During orientation mapping runs, participants detected and reported the presence of a

359     target letter (an X or Y) that appeared at an unpredictable time during the interval from +750

360     msec to +2250 ms following stimulus onset. Responses were made on a USB-compatible

361     number pad. During category discrimination runs, participants ignored the RSVP stream and

362     instead reported whether the orientation of the bar aperture was an exemplar from category "1"

363     or category "2". As in Experiment 1, we randomly designated one of the 15 possible stimulus

364     orientations as the category boundary such that the seven orientations counterclockwise to this

365     value were assigned to Category 1 and the seven orientations clockwise to this value were

366     assigned to Category 2. Participants could respond at any point during the trial, but the stimulus

367     was presented for a total of 3000 msec. Trials were separated by a 2.5 – 3.25 sec inter-trial-

368     interval (randomly selected from a uniform distribution on each trial). Each participant

369     completed four (N = 1), five (N = 10), six (N = 8), seven (N = 8), or eight (N = 1) blocks of the

370     category task and three (N = 1), four (N = 1), five (N = 5), six (N = 12), seven (N = 8), or eight

371     (N = 1) blocks of the orientation mapping task.

372  *EEG Acquisition and Preprocessing.* Participants were seated in a dimly lit, sound-attenuated,

373  and electrically shielded recording chamber (ETS Lindgren) for the duration of the experiment.

374  Continuous EEG was recorded from 128 Ag-AgCl⁻ scalp electrodes via a Biosemi "Active Two"

375  system (Amsterdam, Netherlands). The horizontal electrooculogram (EOG) was recorded from

376  additional electrodes placed near the left and right canthi, and the vertical EOG was recorded

377  from electrodes placed above and below the right eye. Additional electrodes were placed over

378  the left and right mastoids. The horizontal and vertical EOG were recorded from electrodes

379  placed over the left and right canthi and above and below the right eye (respectively). Electrode

380  impedances were kept well below 20 kΩ, and recordings were digitized at 1024 Hz.

381        After testing, the entire EEG time series at each electrode was high- and low-pass filtered

382  (3$^{rd}$ order zero-phase forward and reverse Butterworth) at 0.1 and 50 Hz and re-referenced to the

383  average of the left and right mastoids. Data from both tasks were epoched into intervals spanning

384  -1000 to +4000 msec from stimulus onset; the relatively large pre- and post-stimulus epochs

385  were included to absorb filtering artifacts that could affect later analyses. Trials contaminated by

386  EOG artifacts (horizontal eye movements > 2° and blinks) were identified and excluded from

387  additional analyses. Across participants an average of 5.58% (±1.67%) and 8.74% (±1.84%) of

388  trials from the orientation mapping and category discrimination tasks were discarded

389  (respectively). Finally, noisy channels (those with multiple deflections ≥ 100 μV over the course

390  of the experiment) were visually identified and eliminated (mean number of removed electrodes

391  across participants ±1 S.E.M. = 2.25 ± 0.64).

392        Next, we identified a set of electrodes-of-interest (EOIs) with strong responses at the

393  stimulus' flicker frequency (30 Hz). Data from each task were re-epoched into intervals spanning

394  0 to 3000 msec around stimulus onset and averaged across trials and tasks (i.e., RSVP and

18

395    category discrimination), yielding a *k* electrode by *t* sample data matrix. We computed the

396    evoked power at the stimulus' flicker frequency (30 Hz) by applying a discrete Fourier transform

397    to the average time series at each electrode and selected the 32 electrodes with the highest

398    evoked power at the stimulus' flicker frequency for further analysis. These electrodes were

399    typically distributed over occipitoparietal electrode sites (see Figure 12).

400        To isolate stimulus-specific responses, the epoched timeseries at each electrode was

401    resampled to 256 Hz and then bandpass filtered from 29 to 31 Hz (zero-phase forward and

402    reverse $3^{rd}$ order Butterworth). We next estimated a set of complex Fourier coefficients

403    describing the power and phase of the 30 Hz response by applying a Hilbert transformation to the

404    filtered data. To visualize and quantify orientation-selective signals from frequency-specific

405    responses, we first constructed a complex-valued data set $B_1(t)$ (*m* electrodes x $n_{train}$ trials). We

406    then estimated a complex-valued weight matrix $W(t)$ (*m* channels x *k* filters) using $B_1(t)$ and a

407    basis set of idealized orientation-selective filters $C_1$. Finally, we estimated a complex-valued

408    matrix of channel responses $C_2(t)$ (*m* channels x $n_{test}$ trials) given $W(t)$ and complex-valued test

409    data set $B_2(t)$ (*m* electrodes x $n_{test}$ trials) containing the complex Fourier coefficients measured

410    during the category discrimination task. Trial-by-trial and sample-by-sample response functions

411    were shifted in the same manner described above so that category biases would manifest as a

412    rightward (clockwise) shift towards the center of Category 2. We estimated the evoked (i.e.,

413    phase-locked) power of the response at each filter by computing the squared absolute value of

414    the average complex-valued coefficient for each filter after shifting. Categorical biases were

415    quantified using the same curve fitting analysis described in the main text.

416        To obtain an unbiased estimate of orientation selectivity in each electrode, we ensured

417    that the training data set $B_1(t)$ contained an equal number of trials for each stimulus orientation

418   (0-168° in 12° increments). For each participant, we identified the stimulus orientation θ with the

419   *N* fewest repetitions in the orientation mapping data set after EOG artifact removal. Next, we

420   constructed the training data set $B_1(t)$ by randomly selecting (without replacement) 1:*N* trials for

421   each stimulus orientation. Data from this training set were used to estimate a set of orientation

422   weights for each electrode and these weights were in turn used to estimate a response for each

423   hypothetical orientation channel during the category discrimination task. To ensure that our

424   method generalized across multiple combinations of orientation mapping trials, we repeated this

425   analysis 100 times and averaged the results across permutations.

426   **Experiment 3 - EEG**

427   *Participants*. 8 volunteers recruited from the Florida Atlantic University community completed

428   Experiment 3. All participants self-reported normal or corrected-to-normal visual acuity and

429   gave both written and oral informed consent as required by the local Institutional Review Board.

430   Each participant was tested in a single 2-2.5 hour experimental session (the exact duration varied

431   across participants depending on the amount of time needed to set up and calibrate the EEG

432   equipment).

433   *Behavioral Tasks*. Participants performed six blocks of a spatial recall task followed by multiple

434   blocks of a delayed match-to-category (DMC) task. Both tasks used identical stimulus and

435   display geometry. During the spatial recall task, participants were shown a sample display

436   containing a disc (diameter 2.5° from a viewing distance of 60 cm) rendered in one of 12 polar

437   locations (0° to 330° in 30° increments) along the perimeter of an imaginary circle centered at

438   fixation (radius 7.5°). The sample display was shown for 250 ms and followed by a 1750 ms

439   blank delay. At the end of each trial, participants were shown a mouse cursor and instructed to

440   click on the position of the disc shown in the sample display. Participants were instructed to

441    prioritize accuracy over speed, though a 3000 ms response deadline was imposed. Each trial was

442    followed by a 1500-2200 ms blank interval (randomly sampled from a uniform distribution on

443    each trial). Each block featured 72 trials (six repetitions per stimulus position) and lasted

444    approximately six minutes. EEG data recorded during this task were used to train a position-

445    specific inverted encoding model (see below). Each participant completed six blocks of this task.

446        After completing the spatial recall task, participants performed a delayed match-to-

447    category task. Participants were shown stimuli in the same 12 positions used during the spatial

448    recall task. However, for each participant we defined a category boundary such that half of the

449    possible stimulus positions were assigned membership in Category 1 and the remaining half

450    were assigned membership in Category 2. For example, the category boundary could be set such

451    that positions [315, 345, 15, 45, 75, 105] comprised Category 1 while positions [135, 165, 195,

452    225, 255, 285] comprised Category 2. The location of the category boundary was randomly and

453    independently chosen for each participant and held constant throughout the experiment. At the

454    beginning of each trial, a sample disc appeared in one of the 12 possible stimulus locations for

455    250 ms. After a 1750 ms delay period, a probe disc was presented. The probe could occupy any

456    of the 11 stimulus positions not occupied by the sample, and participants were required to judge

457    whether the position of the probe matched the category of the sample stimulus via keypress.

458    Participants were instructed to prioritize accuracy over speed, but a 3000 ms response limit was

459    imposed. Feedback (correct vs. incorrect) was presented at the end of each trial. Participants

460    completed 5 (N = 1) or 8 (N = 7) blocks of 72 trials.

461        *EEG Acquisition and Preprocessing*. Continuous EEG was recorded from 63 Ag/Ag-Cl⁻

462    scalp electrodes via a Brain Products actiCHamp amplifier. An additional electrode was placed

463    over the right mastoid. Data were recorded with a right mastoid reference and later re-referenced

464    to the algebraic mean of the left and right mastoids (10-20 site TP9 served as the left mastoid

465    reference). The horizontal and vertical electrooculogram (EOG) was recorded from electrodes

466    placed on the left and right canthi and above and below the right eye, respectively. All electrode

467    impedances were kept below 15 kΩ, and recordings were digitized at 1000 Hz. Recorded data

468    were bandpass filtered from 1 to 50 Hz (3[rd] order zero-phase forward and reverse Butterworth

469    filters), epoched from a period spanning -1000 to +3000 ms relative to the start of each trial, and

470    baseline corrected from -250 to 0 ms. Muscle and electrooculogram artifacts were removed from

471    the data using independent components analysis (ICA) as implemented in EEGLAB (Delorme &

472    Makeig, 2004). Reconstructions of stimulus locations were computed from the spatial

473    topography of induced alpha-band (8-12 Hz) power measured across 17 occipitoparietal

474    electrode sites: O1, O2, Oz, PO7, PO3, POz, PO4, PO8, P7, P5, P3, P1, Pz, P2, P4, P6, and P8.

475    *Inverted Encoding Model.* Experiment 3 relied on a fundamentally different signal than

476    Experiment 2 (induced-alpha-band activity vs. evoked 30 Hz power, respectively). Following

477    earlier research (Kok et al., 2017; Ester et al., 2018; Nouri & Ester, 2019), we used a variant of

478    the IEM approach described in Experiment 2 to compute location channel responses. We first

479    isolated alpha-band activity, by bandpass filtering the raw EEG time series at each electrode

480    from 8-12 Hz (zero-phase forward and reverse filters as implemented by EEGLAB's "eegfilt"

481    function), yielding a real-valued signal $f(t)$. The analytic representation of $f(t)$ was obtained by

482    applying a Hilbert transformation:

483

$$z(t) = f(t) + if(t)$$

484

485   where $i = \sqrt{-1}$ and $if(\mathrm{t}) = A(t)e^{i\varphi(t)}$. Induced alpha power was computed by extracting and

486   squaring the instantaneous amplitude $A(t)$ of the analytic signal $z(t)$. We modeled alpha power at

487   each scalp electrode as a weighted sum of 12 location-selective channels, each with an idealized

488   tuning curve (a half-wave rectified cosine raised to the $12^{th}$ power). The maximum response of

489   each channel was normalized to 1, thus units of response are arbitrary. The predicted responses

490   of each channel during each trial were arranged in a $k$ channel by $n$ trials design matrix $C$.

491   Separate design matrices were constructed to track the locations of the blue and red discs across

492   trials (i.e., we reconstructed the locations of the blue and red discs separately, then later sorted

493   these reconstructions according to cue condition). The relationship between the data and the

494   predicted channel responses $C$ is given by a general linear model of the form:

495

$$B = WC + N$$

496

497   where B is a $m$ electrode by $n$ trials training data matrix, W is an $m$ electrode by $k$ channel weight

498   matrix, and $N$ is a matrix of residuals (i.e., noise).

499        To estimate $W$, we constructed a "training" data set containing an equal number of trials

500   from each stimulus location (i.e., 45-360° in 45° steps) condition. We first identified the location

501   $\varphi$ with the fewest $r$ repetitions in the full data set after EOG artifact removal. Next, we

502   constructed a training data set $B_{trn}$ ($m$ electrodes by $n$ trials) and weight matrix $C_{trn}$ ($n$ trials by $k$

503   channels) by randomly selecting (without replacement) $1{:}r$ trials for each of the eight possible

504   stimulus locations (ignoring cue condition; i.e., the training data set contained a mixture of

505   neutral and valid trials). The training data set was used to compute a weight for each channel $C_i$

506   via least-squares estimation:

23

507

$$W_i = B_{trn} C_{trn,i}^T (C_{trn,i} C_{trn,i}^T)^{-1}$$

508

509    where $C_{trn,i}$ is an $n$ trial row vector containing the predicted responses of spatial channel $i$ during

510    each training trial.

511         After estimating the weight matrix $W$, we next estimated a set of spatial filters $V$ that

512    capture the underlying channel responses while accounting for correlated variability between

513    electrode sites (i.e., the noise covariance; Kok et al. 2017):

514

$$V_i = \frac{\Sigma_i^{-1} W_i}{W_i^T \Sigma_i^{-1} W_i}$$

515

516    where $\Sigma_i$ is the regularized noise covariance matrix for channel $i$ and estimated as:

517

$$\sum_i = \frac{1}{n-1} \epsilon_i \epsilon_i^T$$

518

519    where $n$ is the number of training trials and $\varepsilon_i$ is a matrix of residuals:

520

$$\epsilon_i = B_{trn} - W_i C_{trn,i}$$

521

522         Estimates of $\varepsilon_i$ were obtained by regularization-based shrinkage using an analytically

523    determined shrinkage parameter (see Blankertz et al. 2011; Kok et al. 2017). An optimal spatial

524    filter $v_i$ was estimated for each channel $C_i$, yielding an $m$ electrode by $k$ filter matrix $V$. Next, we

24

525    constructed a "test" data set $B_{tst}$ ($m$ electrodes by $n$ trials) containing data from all trials not

526    included in the training data set and estimated trial-by-trial channel responses $C_{tst}$ ($k$ channels x $n$

527    trials) from the filter matrix $V$ and the test data set:

528

$$C_{tst} = V^T B_{tst}$$

529

530        Trial-by-trial channel responses were interpolated to 360°, circularly shifted to a common

531    center (0°, by convention), and sorted by category membership. As in Experiments 1 and 2,

532    reconstructions were shifted and aligned so that any bias would manifest as a shift toward

533    Category B (clockwise). Finally, to ensure internal reliability this entire analysis was repeated 50

534    times, and unique (randomly chosen) subsets of trials were used to define the training and test

535    data sets during each permutation. The results were then averaged across permutations.

536    *Eye Movement Control Analyses – Experiments 2 and 3*. Systematic biases in eye position can

537    contribute to orientation and location performance (e.g., Quax et al., 2019). We did not collect

538    eye position data from Experiment 1 (fMRI). However, different tasks were used to train and test

539    the encoding model, which can be an effective way of mitigating the effects of eye movements

540    on stimulus decoding (Mostert et al., 2018). We also collected electrooculogram (EOG) data

541    during Experiments 2 and 3 (EEG). To examine whether eye position varied as a function of

542    stimulus position during these experiments, we regressed trial-by-trial horizontal EOG

543    recordings (in μV) onto the orientation of a to-be-categorized stimulus (Experiment 2) or the

544    location of a to-be-categorized disc (Experiment 3). In both experiments, we identified and

545    excluded trials contaminated by large horizontal EOG artifacts ($\geq 40$ μV, which corresponds to a

546    horizontal displacement of 2.5° assuming a voltage threshold of 16 μV per degree; Lins et al.,

547     1993), but smaller variations in eye positions – for example, along the inner stimulus aperture –

548     may have escaped detection. Using Experiment 2 as an example, we considered two possibilities.

549     First, participants may have foveated the inner aperture of the stimulus at a polar location

550     matching its orientation. To illustrate, participants could foveate the inner aperture of a 45°

551     stimulus at a polar angle of 45° or 225°; likewise, they could foveate the inner aperture of a 168°

552     stimulus at a polar angle of 168° or 348°. Second, participants may have foveated the inner

553     aperture of each stimulus matching the center of the category it belonged to. We tested these

554     possibilities by calculating predicted horizontal eye positions under the assumption that

555     participants foveated the inner stimulus aperture at locations matching its orientation or the

556     center of the relevant category. Specifically, we converted records of stimulus orientation (or the

557     center of the category to which the stimulus belonged) to polar format and scaled the resulting

558     estimates by the radius of the inner stimulus aperture, then regressed these estimates onto

559     horizontal EOG activity (in μV). If there is a systematic relationship between eye position and

560     either stimulus orientation or category at any point during a trial, then this analysis should yield

561     regression coefficients reliably greater than 0 μV. Identical analyses were used to examine

562     systematic relationships between horizontal eye position and stimulus location in Experiment 3.

563

564

565                                          **Results**

566      *Experiment 1 - fMRI*

567      We trained eight human volunteers to categorize a set of orientations into two groups, Category 1

568      and Category 2. The stimulus space comprised a set of 15 oriented stimuli, spanning 0-168° in

569      12° increments (Figure 1A-B). For each participant, we randomly designated one of these 15

570      orientations as a category boundary such that the seven orientations anticlockwise to the

571      boundary were assigned membership in Category 1 and the seven orientations clockwise to the

572      boundary were assigned membership in Category 2. Each participant completed a one-hour

573      training session prior to scanning. Each participant's category boundary was kept constant across

574      all behavioral training and scanning sessions. Many participants self-reported that they learned

575      the rule delineating the categories in one or two 5-minute blocks of trials. Consequently, task

576      performance measured during scanning was extremely high, with errors and slow responses

577      present only for exemplars immediately adjacent to the category boundary (Figure 1C-D).

578      During each scanning session, participants performed the category discrimination task and an

579      orientation model estimation task where they were required to report the identity of a target letter

580      embedded within a rapid stream presented at fixation while a task-irrelevant grating flickered in

581      the background. Data from this task were used to compute an unbiased estimate of orientation

582      selectivity for each voxel in visual areas V1-hV4v/V3A (see below).

583            We first examined whether category training increased the similarity of activation

584      patterns evoked by exemplars from the same category (i.e., acquired similarity). We tested this

585      by training a linear decoder (support vector machine) to discriminate between activation patterns

586      associated with exemplars at the center of each category (48° from the boundary), then used the

587      trained classifier to predict the category membership of exemplars immediately adjacent to the

588  category boundary (±12°; Figure 2A). This analysis was performed separately for the orientation

589  mapping and category discrimination tasks. We reasoned that if category training homogenizes

590  activation patterns evoked by members of the same category, then decoding performance should

591  be superior during the category discrimination task relative to the orientation mapping task. This

592  is precisely what we observed (Figure 2B). For example, near-boundary decoding performance

593  in V1 was reliably above chance during the category discrimination task ($p < 0.0001$, false-

594  discovery-rate-corrected bootstrap test), but not during the orientation mapping task when the

595  category boundary was irrelevant and the oriented stimulus was unattended ($p = 0.38$).

596  Importantly, the absence of robust decoding performance during the orientation mapping task

597  cannot be attributed to poor signal, as a decoder trained and tested on activation patterns

598  associated with exemplars at the center of each category (Figure 2C) yielded above-chance

599  decoding during both behavioral tasks (Figure 2D; $M = 0.58$ and $0.69$ for the mapping and

600  discrimination tasks, respectively; $p < 0.01$, bootstrap test). Collectively, these results suggest

601  that category training can alter population-level responses at very early stages of the visual

602  processing hierarchy.

603      To better understand how category training influences orientation-selective activation

604  patterns in early visual cortical areas, we used an inverted encoding model (Brouwer & Heeger,

605  2011) to generate model-based reconstructed representations of stimulus orientation from these

606  patterns. For each visual area (e.g., V1), we first modelled voxel-wise responses measured during

607  the orientation mapping task as a weighted sum of idealized orientation channels, yielding a set

608  of weights that characterize the orientation selectivity of each voxel (Figure 3A). In the second

609  phase of the analysis, we reconstructed trial-by-trial representations of stimulus orientation by

610  combining these weights with the observed pattern of activation across voxels measured during

611     each trial of the category discrimination task, resulting in single-trial reconstructed channel

612     response function that contains a representation of stimulus orientation for each ROI on each trial

613     (Figure 3B). Finally, we sorted trial-by-trial reconstructions according to category membership

614     such that any bias would manifest as a clockwise (rightward) shift of the orientation

615     representation towards the center of Category 2 and quantified biases towards this category using

616     a curve-fitting analysis (Methods).

617          Note that stimulus orientation was irrelevant during the orientation mapping task used for

618     model weight estimation. We therefore reasoned that voxel-by-voxel responses evoked by each

619     oriented stimulus would be largely uncontaminated by its category membership. Indeed, the

620     logic of our analytical approach rests on the assumption that orientation-selective responses are

621     quantitatively different during the orientation mapping and category discrimination tasks: if

622     identical category biases are present in both tasks then the orientation weights computed using

623     data from either task will capture that bias and reconstructed representations of orientation will

624     not exhibit any category shift. This is precisely what we observed when we used a cross-

625     validation approach to reconstruct stimulus representations separately for the orientation

626     mapping and category discrimination tasks (Figure 4).

627          As shown in Figure 5, reconstructed representations of orientation in visual areas V1, V2,

628     and V3 were systematically biased away from physical stimulus orientation and towards the

629     center of the appropriate category (shifts of 22.13°, 26.65°, and 34.57°, respectively; $P < 0.05$,

630     bootstrap test, false-discovery-rate [FDR] corrected for multiple comparisons across regions; see

631     Figure 6 for separate reconstructions of Category 1 and Category 2 orientations and Figure 7 for

632     participant-by-participant reconstructions plotted by visual area). Similar, though less robust

633     biases were also evident in hV4v and V3A (mean shifts of 9.73° and 6.45°, respectively; $p >$

29

634  0.19). A logistic regression analysis established that categorical biases in V1-V3 were strongly

635  correlated with variability in overt category judgments (Figure 8). That is, trial-by-trial category

636  judgments were more strongly associated with the responses of orientation channels near the

637  center of each category rather than those near the physical orientation of the stimulus.

638  Importantly, because the location of the boundary separating categories 1 and 2 was randomly

639  selected for each participant, it is unlikely that categorical biases shown in Figure 5 reflect

640  intrinsic biases in stimulus selectivity in early visual areas (e.g., due to oblique effects; Sun et al.,

641  2013).

642      The category biases shown in Figure 5 may be the result of an adaptive process that

643  facilitates task performance by enhancing the discriminability of physically similar but

644  categorically distinct stimuli. Consider a hypothetical example where an observer is tasked with

645  discriminating between two physically similar exemplars on opposite sides of a category

646  boundary. Discriminating between these alternatives should be challenging as each exemplar

647  evokes a similar and highly overlapping response pattern. However, discrimination performance

648  could be improved if the responses associated with each exemplar are made more separable via

649  acquired distinctiveness following training (or equivalently, an acquired similarity between

650  exemplars adjacent to the category boundary and exemplars near the center of each category).

651  Similar changes would be less helpful when an observer is tasked with discriminating between

652  physically and categorically distinct exemplars, as each exemplar already evokes a dissimilar and

653  non-overlapping response. From these examples, a simple prediction can be derived: categorical

654  biases in reconstructed representations of orientation should be largest when participants are

655  shown exemplars adjacent to the category boundary and progressively weaker when participants

656  are shown exemplars further away from the category boundary.

657        We tested this possibility by sorting stimulus reconstructions according to the angular

658 distance between stimulus orientation and the category boundary (Figure 9). As predicted,

659 reconstructed representations of orientations adjacent to the category boundary were strongly

660 biased by category membership, with larger biases for exemplars nearest to the category

661 boundary ($\mu$ = 42.62°, 24.16°, and 20.12° for exemplars located 12°, 24°, and 36° from the

662 category boundary, respectively; FDR-corrected bootstrap $p < 0.0015$), while reconstructed

663 representations of orientations at the center of each category exhibited no signs of bias ($\mu$ = -

664 3.98°, $p = 0.79$; the direct comparison of biases for exemplars adjacent to the category boundary

665 and in the center of each category was also significant; $p < 0.01$). Moreover, the relationship

666 between average category bias and distance from the category boundary was well-approximated

667 by a linear trend (slope = -14.38°/step; $r^2 = 0.96$). Thus, category biases in reconstructed

668 representation are largest under conditions where they would facilitate behavioral performance

669 and absent under conditions where they would not.

670        Category-selective signals have been identified in multiple brain areas, including portions

671 of lateral occipital cortex, inferotemporal cortex, posterior parietal cortex, and lateral prefrontal

672 cortex (Sigala & Logothetis, 2002; Freedman et al., 2011; Freedman & Assad, 2006; Folstein et

673 al., 2012; Davis & Poldrack, 2013; Pourtois et al., 2008; Mack et al., 2013). We identified

674 category selective information in many of these same regions using a whole-brain searchlight-

675 based decoding analysis where a classifier was trained to discriminate between exemplars from

676 Category 1 and Category 2 (independently of stimulus orientation; Figure 10 and Methods).

677 Next, we used the same inverted encoding model described above to reconstruct representations

678 of stimulus orientation from activation patterns measured in each area during each of the

679 orientation mapping and category discrimination tasks (reconstructions were computed using a

680     "leave-one-participant-out" cross-validation routine to ensure that reconstructions were

681     independent of the decoding analysis used to define category-selective ROIs). We were able to

682     reconstruct representations of stimulus orientation in many of these regions during the category

683     discrimination task, but not during the orientation mapping task (where stimulus orientation was

684     task-irrelevant; Figure 11). This is perhaps unsurprising as representations in many mid-to-high

685     order cortical areas are strongly task-dependent (e.g., Silver et al., 2005).  As our analytical

686     approach requires an independent and unbiased estimate of each voxel's orientation selectivity

687     (e.g., during the orientation mapping task), this meant that we were unable to probe categorical

688     biases in reconstructed representations in these regions.

689     *Experiment 2 - EEG*

690             Due to the sluggish nature of the hemodynamic response, the category biases shown in

691     Figures 5 and 9 could reflect processes related to decision making or response selection rather

692     than stimulus processing. In a second experiment, we evaluated the temporal dynamics of

693     category biases using EEG. Specifically, we reasoned that if the biases shown in Figures 5 and 9

694     reflect processes related to decision making, response selection, or motor planning, then these

695     biases should manifest only during a period shortly before the participants' response.

696     Conversely, if the biases are due to changes in how sensory neural populations encode features,

697     they should be evident during the early portion of each trial. To evaluate these alternatives, we

698     recorded EEG while a new group of 28 volunteers performed variants of the orientation mapping

699     and categorization tasks used in the fMRI experiment. On each trial, participants were shown a

700     large annulus of iso-oriented bars that flickered at 30 Hz (i.e., 16.67 ms on, 16.67 ms off; Figure

701     12A). During the orientation mapping task, participants detected and reported the identity of a

702     target letter (an X or a Y) that appeared in a rapid series of letters over the fixation point.

703    Identical displays were used during the category discrimination task, with the caveat that

704    participants were asked to report the category of the oriented stimulus while ignoring the letter

705    stream.

706        The 30 Hz flicker of the oriented stimulus elicits a standing wave of frequency-specific

707    sensory activity known as a steady-state visually-evoked potential (SSVEP, Vialatte et al., 2010;

708    Figure 12B). The coarse spatial resolution of EEG precludes precise statements about the cortical

709    source(s) of these signals (e.g., V1, V2, etc.). However, to focus on visual areas (rather than

710    parietal or frontal areas) we deliberately entrained stimulus-locked activity at a relatively high

711    frequency (30 Hz). Our approach was based on the logic that coupled oscillators can only be

712    entrained at high frequencies within small local networks, while larger or more distributed

713    networks can only be entrained at lower frequencies due to conduction delays (Breakspear et al.,

714    2010). Indeed, a topographic analysis showed that evoked 30 Hz activity was strongest over a

715    localized region of occipitoparietal electrode sites. (Figure 12C). As in Experiment 1,

716    participants learned to categorize stimuli with a high degree of accuracy, with errors and slow

717    responses present only for exemplars adjacent to a category boundary (Figure 12D-E)

718        We computed the power and phase of the 30 Hz SSVEP response across each 3,000 ms

719    trial and then used these values to reconstruct a time-resolved representation of stimulus

720    orientation (Garcia et al., 2013). Our analysis procedure followed that used in Experiment 1: In

721    the first phase of the analysis, we rank-ordered scalp electrodes by 30 Hz power (based on a

722    discrete Fourier transform spanning the 3000 ms trial epoch, averaged across all trials of both the

723    orientation mapping and category discrimination tasks). Responses measured during the

724    orientation mapping task were used to estimate a set of orientation weights for the 32 electrodes

725    with the strongest SSVEP signals (i.e., those with the highest 30 Hz power; see Figure 12C) at

33

726    each timepoint. In the second phase of the analysis, we used these timepoint-specific weights and

727    corresponding responses measured during each trial of the category discrimination task across all

728    electrodes to compute a time-resolved representation of stimulus orientation (Figure 13A-B). We

729    reasoned that if the categorical biases shown in Figures 5 and 9 reflect processes related to

730    decision making or response selection, then they should emerge gradually over the course of

731    each trial. Conversely, if the categorical biases reflect changes in sensory processing, then they

732    should manifest shortly after stimulus onset. To test this possibility, we computed a temporally

733    averaged stimulus reconstruction over an interval spanning 0 to 250 ms after stimulus onset

734    (Figure 14B). A robust category bias was observed (M = 23.35°; p = 0.014; bootstrap test)

735    suggesting that the intent to categorize a stimulus modulates how neural populations in early

736    visual areas respond to incoming sensory signals.

737        Importantly, the bandpass filter used to isolate 30 Hz activity will distort temporal

738    characteristics of the raw EEG signal by some amount. We estimated the extent of this distortion

739    by generating a 3 second, 30 Hz sinusoid with unit amplitude (plus 1 second of pre-and post-

740    signal zero padding) and running it through the same filters used in our analysis path. We then

741    computed the time at which the filtered signal reach 25% of maximum. For an instantaneous

742    filter, this should occur at exactly 1 second (due to the pre- and post-signal zero-padding). We

743    estimated a signal onset of ~877 ms, or 123 ms prior to the start of the signal. Therefore, if

744    reconstruction amplitude is greater than zero at time t, then we can conclude that the pattern of

745    scalp activity used to generate the stimulus reconstruction contained reliable orientation

746    information at time $t \pm 125$ ms. The same logic applies to estimates of reconstruction bias as the

747    reconstructions are based on data filtered using the same parameters. Importantly, we also

748    verified that there was no categorical bias in stimulus reconstructions prior to stimulus onset

749    (Figure 14), nor were categorical biases present when we reconstructed stimulus representations

750    using data from the orientation mapping and category discrimination tasks separately (Figure

751    15).

752    *Ruling out contributions from eye movements.* We identified and removed trials contaminated by

753    large EOG artifacts (blinks and eye movements greater than ~2°). However, small and consistent

754    eye movement patterns could nevertheless contribute to the orientation reconstructions reported

755    here. We examined this possibility by testing whether participants foveated the inner aperture of

756    the stimulus at polar locations matching its orientation (Figure 16A) or at polar locations

757    matching the center of the appropriate category (A vs B; Figure 16B; see Methods for details).

758    No systematic differences in eye position were observed as a function of stimulus orientation or

759    category membership (Figure 16), suggesting that eye movements were not a major contributor

760    to orientation-specific reconstructions.

761    *Experiment 3 - EEG*

762         The results of Experiments 1 and 2 suggest that category learning can bias stimulus-

763    specific representations encoded by occipitoparietal cortical areas. However, an alternative

764    explanation is that the biases shown in Figures 5, 9, and 13 reflect mechanisms of response

765    selection or decision making independent of categorical processing. Experiment 3 examined this

766    possibility by examining categorical biases in stimulus-specific memory representations while

767    participants performed a delayed match-to-category (DMC) task. A schematic of the task is

768    shown in Figure 17A-B. At the beginning of each trial a sample disc rendered in one of 12

769    possible stimulus locations (15-345° polar angle in 30° along the perimeter of an imaginary

770    circle). Half of the disc positions were assigned membership in Category 1, while the remaining

771    half of disc positions were assigned membership in Category 2 (Figure 17A). Participants

772     remembered the position of the sample disc over a blank delay, then judged whether a probe disc

773     was rendered in a position matching the category of the sample disc. The location of the category

774     boundary was randomly determined for each participant, and response feedback (correct vs.

775     incorrect) was provided after every trial. Like Experiment 2, participants were not trained on the

776     DMC task prior to testing and learned to associate specific positions with specific categories

777     through feedback. Before completing the DMC task, participants also completed a spatial

778     working memory task. Display and stimulus geometry were identical during the spatial memory

779     task and the DMC task. On each trial a sample disc was rendered in one of the same 12 positions

780     used during the DMC task. After a short delay, participants recalled the location of the sample

781     disc via mouse click.

782         Following earlier work (e.g., Foster et al., 2016; Samaha et al., 2016; Ester et al., 2018;

783     Nouri & Ester, 2019), we used spatiotemporal patterns of induced alpha-band (8-12 Hz) activity

784     over occipitoparietal electrode sites to track the contents of spatial working memory during the

785     recall and DMC tasks. Specifically, we modeled sample-by-sample estimates of alpha band

786     activity recorded during the spatial recall task as a combination of 12 location filters, each with

787     an idealized tuning curve (a cosine raised to the $12^{th}$ power). The result of this step is a set of

788     weights that characterizes the location preferences of each scalp electrode. Next, we used these

789     weights and spatiotemporal patterns of alpha-band activity recorded during the DMC task to

790     compute an expected response for each location filter, yielding a time-resolved estimate of

791     stimulus position. Trial-by-trial response functions were shifted to a common center (0° by

792     convention), averaged, and arranged such that any category bias would manifest as a clockwise

793     or positive shift towards the center of Category 2.

794    As expected, a robust category bias was observed during the delay period of the DMC

795    task (Figure 17C), though unlike Experiment 2 the bias seemed to emerge gradually over the

796    course of the delay period. To quantify this bias, we averaged channel responses from period

797    0.25 to 2.0 sec after onset of the sample display and fit the resulting function with an

798    exponentiated cosine (*Quantification of Bias in Orientation Representations*, Methods). Mean

799    reconstruction centers were reliably greater than 0° (M = 10.55°; $p$ = 0.002, bootstrap test),

800    indicating a robust bias towards the center of the relevant category. Importantly, this bias cannot

801    be explained by mechanisms associated with decision making and response selection:

802    participants could not plan or implement a response until the probe stimulus was presented at the

803    end of the delay period. This result further suggests that the results of Experiments 1 and 2

804    cannot be wholly explained by mechanisms of response selection or bias.

805    *Assessing contributions from eye movements.* We identified and removed electrooculogram

806    artifacts from the data via independent components analysis. However, small and consistent eye

807    movement patterns opaque to ICA could nevertheless contribute to the location reconstructions

808    reported here. We examined this possibility by regressing time-resolved estimates of horizontal

809    EOG activity onto remembered stimulus locations. As shown in Figure 18, the regression

810    coefficients linking eye position with remembered locations were indistinguishable from 0 for

811    the duration of each trial, suggesting that eye movements were not a major determinant of

812    location reconstructions.

**Discussion**

813

814      Our findings suggest that category learning shapes information processing at the earliest

815 stages of the visual system. The results of Experiment 1 showed that representations of a to-be-

816 categorized stimulus encoded by population-level activity in early visual cortical areas were

817 systematically biased by their category membership. These biases were correlated with overt

818 category judgments and were largest for exemplars adjacent to the category boundary. The

819 results of Experiments 2 and 3 demonstrate that similar biases are present in orientation- and

820 location-specific reconstructions computed by human scalp EEG data, and further suggest that

821 our findings cannot be explained by response bias, motor planning, or eye movements.

822      The categorical biases reported here are strongly task dependent, and therefore must

823 reflect changes in responses caused by transient top-down factors rather than long-term changes

824 in feature or location selectivity. However, the effects of these top down modulations are

825 fundamentally different from task-dependent modulations reported elsewhere. In one example,

826 Ester et al. (2016) asked participants to attend the orientation or luminance of a peripheral

827 grating and found both multiplicative and additive enhancements of orientation-specific

828 reconstructions during the attend orientation condition relative to the attend luminance condition,

829 but no evidence for a shift like the one reported here. In a different study, Byers and Serences

830 (2014) examined changes in orientation-specific reconstructions before and after participants

831 underwent extensive training (10 1-hour sessions) in a challenging orientation discrimination

832 task. We observed changes in the amplitude (i.e., signal-to-noise ratio) of orientation-specific

833 reconstructions following training, but no evidence for a shift like the one reported in the current

834 study. In other studies, Scolari et al. (2012; 2014) examined changes in orientation-specific

835 reconstructions when participants performed fine-grained and coarse-grained orientation

836  discrimination tasks. Participants viewed two oriented gratings in sequence and judged whether

837  they were identical. During one experiment participants were cued to how the second grating

838  might differ from the first (clockwise vs. counterclockwise rotation), while in a second

839  experiment they were not. During the fine-grained discrimination task, the authors observed a

840  modest shift in orientation-specific reconstructions towards "off-target" neural populations that

841  maximally discriminated between two oriented stimuli, but only when participants were cued to

842  expect a clockwise or counterclockwise rotation. While this type of modulation is desirable while

843  performing a fine-discrimination task, it is qualitatively different than the shifts we report in the

844  current experiment, as participants have no way of anticipating what orientation will be

845  presented on each trial, nor the difference between that orientation and the category boundary.

846  Moreover, the shifts reported by Scolari et al. (2012) during fine discriminations were relatively

847  modest – at most few degrees. We report an opposite pattern of findings, where shifts are largest

848  for oriented exemplars immediately adjacent to the category boundary. Thus, while other studies

849  have documented task-dependent changes in orientation-specific reconstructions, those studies

850  have failed to reveal shifts in reconstructed representations (Ester et al., 2016; Byers & Serences

851  2014) or have revealed modest shifts that follow different patterns from those reported here

852  (Scolari et al. 2012).

853      Several mechanisms may be responsible for our findings. One possibility is that the

854  orientation preferences of single-units (or populations of units) are systematically shifted towards

855  the center of each category during the category discrimination task, much in the same way that

856  neurons in the rodent auditory system exhibit emergent selectivity for categorically different

857  stimuli (e.g., Xin et al., 2019)  or in the same way that the spectral preferences of neural

858  populations are biased by feature-based attention (David et al., 2008; Cukur et al., 2012). These

39

859   shifts are relatively small at the single unit level but could be amplified by a read-out

860   mechanisms that integrate the responses of large neural populations. A second possibility is that

861   participants strategically apply gain to neural populations that maximally discriminate between

862   to-be-categorized exemplars during the category discrimination task. Here there are no changes

863   in the spectral preferences of single units, but the responses of neurons that respond to stimuli

864   near the category boundary are amplified. These alternatives are not mutually exclusive; nor is

865   this an exhaustive list. Our data cannot resolve these possibilities. For example, several different

866   patterns of single-unit gain changes and/or tuning shifts can produce identical responses in a

867   single fMRI voxel, and different patterns of single-voxel modulation could produce categorical

868   biases in multivariate stimulus reconstructions (see, e.g., Sprague et al., 2018 for a detailed

869   discussion of this issue). Ultimately, targeted experiments that combine non-invasive

870   measurements of brain activity with careful psychophysical measurements and detailed model

871   simulations will be needed to conclusively identify the mechanisms responsible for the category

872   biases we have reported here.

873       Our findings appear to conflict with results from nonhuman primate research which

874   suggests that sensory cortical areas do not encode categorical information. However, there is

875   reason to suspect that mechanisms of category learning might be qualitatively different in human

876   and non-human primates. For example, our participants learned to categorize stimuli based on

877   performance feedback after approximately 10 minutes of training. In contrast, macaque monkeys

878   typically require six months or more of training using a similar feedback scheme to reach a

879   similar level of performance, and this extensive amount of training may influence how neural

880   circuits code information (e.g., Itthipurripat et al., 2017; Birman & Gardner, 2015). Moreover,

881   there is growing recognition that the contribution(s) of sensory cortical areas to performance on a

882 visual task are highly susceptible to recent history and training effects (Itthipurripiat et al., 2017,

883 Chen et al., 2016; Liu & Pack, 2017). In one example (Liu & Pack, 2017), extensive training was

884 associated with a functional substitution of human visual area V3a for MT+ in discriminating

885 noisy motion patches. Thus, training effects may help explain why previous electrophysiological

886 experiments have found category-selective responses in association but not sensory cortical

887 areas.

888 Studies of categorization in non-human primates have typically employed variants of a

889 delayed match to category task, where monkeys are shown a sequence of two exemplars

890 separated by a blank delay interval and asked to report whether the category of the second

891 exemplar matches the category of the first exemplar. The advantage of this task is that it allows

892 experimenters to decouple category-selective signals from activity related to decision making,

893 response preparation, and response execution. However, this same advantage also precludes

894 examinations of whether and/or how top-down category-selective signals interact with bottom-up

895 stimulus-specific signals. We made no effort to decouple category-selective and decision-related

896 signals in Experiments 1-2, and thus the category biases observed in those studies could reflect

897 mechanisms of decision making, response selection, or motor planning. The results of

898 Experiment 3 conflict with this interpretation by demonstrating that robust category biases are

899 present during the memory period of a delayed match-to-category task (Freedman & Assad,

900 2006).

901 Previous studies have identified cortical modules selective for faces (Kanwisher et al.,

902 1997), locations (Epstein & Kanwisher, 1998), actions (Astafiev et al., 2004; Huth et al., 2012),

903 bodies (Downing et al., 2001); animacy (Konkle & Caramazza, 2013) and size (Konkle &

904 Caramazza, 2013). Other category distinctions (e.g., tools vs. cars) lack specialized processing

905　modules but can be decoded from multivoxel patterns in multiple occipitotemporal regions (e.g.,

906　Folstein et al., 2012). Our findings complement these studies by demonstrating that learning a

907　novel and arbitrary category rule is correlated with rapid and reversible changes in stimulus-

908　specific information processing at even earlier stages of the cortical visual processing hierarchy,

909　including V1 (see also Brouwer & Heeger, 2009; 2013). Category-dependent changes in early

910　visual areas may contribute to more complex forms of category selectivity exhibited by upstream

911　cortical areas, including portions of lateral occipital and inferotemporal cortex. This possibility

912　awaits further scrutiny.

913　　　　To summarize, we have shown that learning a novel and arbitrary category rule based on

914　a simple visual feature (orientation or location) correlates with rapid and reversible changes in

915　sensory and mnemonic representations encoded by regions in early occipitoparietal cortex. These

916　changes correlate with participants' overt category judgments, are largest for exemplars adjacent

917　to a category boundary, and cannot be explained by decision making or motor preparation.

918　Collectively, these results provide novel and important evidence suggesting that category

919　learning induces rapid-yet-reversable changes in information processing at early stages of the

920　cortical visual processing hierarchy.

921

922 **References**

923 Ashby FG, Maddox WT (2005) Human category learning. *Annu Rev Psychol* 56:148-178.

924 Astafiev SV, Stanley CM, Shulman GL, Corbetta M (2004) Extrastriate body area in human

925    occipital cortex responds to the performance of motor actions. *Nat Neurosci* 7:542-548

926 Birman D, Gardiner JL (2015) Parietal and prefrontal: categorical differences? *Nat Neurosci*

927    19:5-7

928 Breakspear M, Heitmann S, Daffertshofer A (2010) Generative models of cortical oscillations:

929    Neurobiological implication of the Kuramoto model. *Front Hum Neurosci* 4:190

930 Brouwer GJ, Heeger DJ (2011) Decoding and reconstructing color from responses in human

931    visual cortex. *J Neurosci* 29:13992-14003

932 Brouwer GJ, Heeger DJ (2011) Cross-orientation suppression in human visual cortex. *J*

933    *Neurophysiol* 106:2108-2119.

934 Brouwer GJ, Heeger DJ (2013) Categorical clustering of the neural representation of color. *J*

935    *Neurosci* 33:15454-15465

936 Byers A, Serences JT (2014) Enhanced attentional gain as a mechanism for generalized

937    perceptual learning in human visual cortex. *J Neurophysiol* 112:1217-1227

938 Chang C-C, Lin CJ (2011) LIBSVM: A library for support vector machines. *ACM Trans Intell*

939    *Syst Technol* 2(27):1-27

940 Chen N, Cai P, Zhou T, Thompson B, Fang F (2016) Perceptual learning modifies the functional

941    specializations of visual cortical areas. *Proc Natl Acad Sci USA* 113:5724-5729

942 Cousineau D (2005) Confidence intervals in within-subject designs: A simpler solution to Loftus

943    & Masson's method. *Quant Meth Psych* 1:42-45

944  Cukur T, Nishimoto S, Huth AG, Gallant JL (2013) Attention during natural vision warps

945      semantic representation across the human brain. *Nat Neurosci* 16:763-770

946  David SV, Hayden BY, Mazer JA, Gallant JL (2008) Attention to stimulus features shifts

947      spectral tuning of V4 neurons during natural vision. *Neuron* 59:509-521

948  Davis T, Poldrack RA (2013) Quantifying the internal structure of categories using a neural

949      typicality measure. *Cereb Cortex* 24:1720-1737.

950  Downing PE, Jiang Y, Shuman M, Kanwisher N (2001) A cortical area selective for visual

951      processing of the human body. *Science* 293:2470-2473

952  Epstein R, Kanwisher N (1998) A cortical representation of the local visual environment. *Nature*

953      392:598

954  Ester EF, Sprague TC, Serences JT (2015) Parietal and frontal cortex encode stimulus-specific

955      mnemonic representations during visual working memory. *Neuron* 87:893-905.

956  Ester EF, Nouri A, Rodriguez L (2018) Retrospective cues mitigate information loss in human

957      cortex during working memory storage. *J Neurosci* 38:8538-8548

958  Esterman M, Tamber-Rosenau BJ, Chiu Y-C, Yantis S (2010) Avoiding non-independence in

959      fMRI data analysis: Leave one subject out. *NeuroImage* 50:572-576

960  Folstein JR, Palmeri TJ, Gauthier I (2012) Category learning increases discriminability of

961      relevant object dimensions in visual cortex. *Cereb Cortex* 23:714-823.

962  Freedman DJ, Riesenhuber M, Poggio T, Miller EK (2011) Categorical representation of visual

963      stimuli in the primate prefrontal cortex. *Science* 291:312-316

964  Freedman DJ, Assad JA (2006) Experience-dependent representation of visual categories in

965      parietal cortex. *Nature* 443:85-88.

966  Friston K (2010) The free-energy principle: A unified brain theory? *Nat Rev Neurosci* 11:127-

967      138

968  Garcia JO, Sreenivasan R, Serences JT (2013) Near-real-time feature-selective modulations in

969      human cortex. *Curr Biol* 23:515-522

970  Gardiner JL, Liu T (2019) Inverted encoding models reconstruct an arbitrary model response, not

971      the stimulus. *eNeuro*

972  Goldstone RL (1994) Influence of categorization on perceptual discrimination. *J Exp Psychol*

973      *Gen* 123:178-200

974  Goldstone RL (1998) Perceptual Learning. *Annu Rev Psychol* 49:585-612

975  Huth AG, Nishimoto S, Vu AT, Gallant JL (2012) A continuous semantic space describes the

976      representation of thousands of object and action categories across the human brain.

977      *Neuron* 76:1210-1224

978  Itthipurripat S, Cha K, Byers A, Serences JT (2017) Two different mechanisms support selective

979      attention at different phases of training. *PLOS Biology*

980  Jazayeri M, Movshon JA (2006) Optimal representation of sensory information by neural

981      populations. *Nat Neurosci* 9:690-696

982  Kanwisher N, McDermott J, Chun MM (1997) The fusiform face area: a module in human

983      extrastriate cortex specialized for face perception. *J Neurosci* 17:4302-4311

984  Konkle T, Caramazza A (2013) Tripartite organization of the ventral stream by animacy and

985      object size. *J Neurosci* 33:10235-10242

986  Kleiner M, Brainard D, Pelli D (2007) What's new in Psychtoolbox-3. *Perception* 36(14)

987  Koida K, Komatsu H (2007) Effects of task demands on the responses of color-selective neurons

988      in the inferior temporal cortex. *Nat Neurosci* 10:108-116

989   Liu LD, Pack CC (2017) The contribution of area MT to visual motion perception depends on

990       training. *Neuron* 95:436-446

991   Livingston K, Andrews J, Harnad S (1998) Categorical perception effects induced by category

992       learning. *J Exp Psychol Learn Mem Cogn* 24:732-753.

993   Mack ML, Preston AR, Love BC (2013) Decoding the brain's algorithm for categorization from

994       it's neural implementation. *Curr Biol* 23:2023-2027.

995   Martinez-Trujillo JC, Treue S (2004) Feature-based attention increases the selectivity of

996       population responses in primate visual cortex. *Curr Biol* 14:744-751.

997   Navalpakkam V, Itti L (2007) Search goal tunes visual features optimally. *Neuron* 53:605-617

998   Newell FN, Bulthoff HH (2002) Categorical perception of familiar objects. *Cognition* 85:113-

999       143

1000  Nouri A, Ester EF (2019) Recovery of information from latent memory stores decreases over

1001      time. *Cogn Neurosci* doi: 10.1080/17588928.2019.1617258

1002  Pourtois G, Schwartz S, Spiridon M, Martuzzi R, Vuilleumier P (2008) Object representations

1003      for multiple visual categories overlap in lateral occipital and medial fusiform cortex.

1004      *Cereb Cortex* 19:1806-1819

1005  Rao RPN, Ballard DH (1999) Predictive coding in the visual cortex: A functional interpretation

1006      of some extra-classical receptive-field effects. *Nat Neurosci* 2:79-87

1007  Scolari M, Serences JT (2010) Basing perceptual decision on the most informative sensory

1008      neurons. *J Neurophysiol* 104:2266-2273

1009  Scolari M, Byers A, Serences JT (2012) Optimal deployment of attentional gain during fine

1010      discriminations. *J Neurosci* 32:7723-7733

1011    Sigala N, Logothetis NK (2002) Visual categorization shapes feature selective in the primate

1012         temporal cortex. *Nature* 415:318-320.

1013    Silver MA, Ress D, Heeger DJ (2005) Topographic maps of visual spatial attention in human

1014         parietal cortex. *J Neurophysiol* 94:1358-1371

1015    Sprague TC, Adam KCS, Foster JJ, Rahmati M, Sutterer DW, Vo VA (2018) Inverted encoding

1016         models assay population-level stimulus representations, not single-unit neural tuning.

1017         *eNeuro* 5(3)

1018    Sun P, Gardiner JL, Costagli M, Ueno K, Waggoner RA, *et al*. Demonstration of tuning to

1019         stimulus orientation in the human cortex: A high-resolution fMRI study with a novel

1020         continuous stimulation paradigm. *Cereb Cortex* 23:1618-1629

1021    Vialatter F-B, Maurice M, Dauwels J, Cichocki A (2010) Steady-state visually evoked

1022         potentials: Focus on essential paradigms and future perspectives. *Prog Neurobiol* 90:418-

1023         438

1024

**Figure 1. Overview of Experiment 1.** (A) Participants viewed displays containing a circular aperture of iso-oriented bars. On each trial, the bars were assigned one of 15 unique orientations from 0-168°. (B) We randomly selected and designated one stimulus orientation as a category boundary (black dashed line), such that the seven orientations counterclockwise from this value were assigned to Category 1 (red lines) and the seven orientations clockwise from this value were assigned to Category 2 (blue lines). (C) After training, participants rarely miscategorized orientations. (D) Response latencies are significantly longer for oriented exemplars near the category boundary (RT = response time; shaded regions in C-D are ±1 within-participant S.E.M.).

**Figure 2. Category Decoding Performance.** (A) We trained classifiers on activation patterns evoked by exemplars at the center of each category boundary during the orientation mapping and category discrimination task (blue lines), then used the trained classifier to predict the category membership of exemplars adjacent to the category boundary (red lines). (B) Decoding accuracy was significantly higher during the category discrimination task relative to the orientation mapping task (p = 0.01), suggesting that activation patterns evoked by exemplars adjacent to the category boundary became more similar to activation patterns evoked by exemplars at the center of each category during the categorization task. The absence of robust decoding performance during the orientation mapping task cannot be attributed to poor signal or a uniform enhancement of orientation representations by attention, as a decoder trained and tested on activation patterns associated with exemplars at the center of each category (C) yielded above-chance decoding during both behavioral tasks (D). Decoding performance was computed from activation patterns in V1. Error bars depict ±1 S.E.M.

**Figure 3. Inverted Encoding Model.** (A) In the first phase of the analysis, we estimated an orientation selectivity profile for each voxel retinotopically organized V1-hV4/V3a using data from an independent orientation mapping task. Specifically, we modeled the response of each voxel as a set of 15 hypothetical orientation channels, each with an idealized response function. (B) In the second phase of the analysis, we computed the response of each orientation channel from the estimated orientation weights and the pattern of responses across voxels measured during each trial of the category discrimination task. The resulting reconstructed channel response function (CRF) contains a representation of the stimulus orientation (example; 24 deg), which we quantified via a curve-fitting procedure.

**Figure 4**. **Reconstructions of stimulus orientation during the orientation mapping task (blue) and the category discrimination task (red)**. Reconstructions were computed using a leave-one-run-out cross validation approach where data from N-1 runs were used to estimate channel weights and data from the remaining run were used to estimate channel responses. This procedure was iterated until all runs had been used to estimate channel responses and the results were averaged over permutations. No categorical biases were observed in any visual area for either task. Shaded regions depict ±1 within-participant S.E.M. a.u., arbitrary units.

1074



1075
1076 **Figure 5. Reconstructed representations of Orientation in Early Visual Cortex**. The vertical
1077 bar at 0° indicates the actual stimulus orientation presented on each trial. Channel response
1078 functions (CRFs) from Category 1 and Category 2 trials have been arranged and averaged such
1079 that any categorical bias would manifest as a clockwise (rightward) shift in the orientation
1080 representation towards the center of Category B. Shaded regions are ±1 within-participant S.E.M
1081 (see Methods). Note change in scale between visual areas V1-V3 and hV4-V3A. a.u., arbitrary
1082 units.
1083

1084
1085
1086 **Figure 6. Stimulus Reconstructions during Category 1 and Category 2 trials.** Shaded regions
1087 are ±1 within-participant S.E.M. a.u., arbitrary units.
1088

**Figure 7. Participant-level Stimulus Reconstructions.** Each panel plots a reconstructed representation of stimulus orientation for a given participant (columns) and visual area (rows). Dashed blue lines are the estimated peak of each reconstruction (obtained via curve-fitting). Ordinate units are arbitrary.

**Figure 8. Categorical Biases predict Choice Behavior.** Each plot shows a logistic regression of each orientation channel's response onto trial-by-trial variability in category judgments. A positive coefficient indicates a positive relationship between an orientation channel's response and the correct category judgment (i.e., Category B), while a negative coefficient indicates a negative relationship between an orientation channel's response and correct category judgment (i.e., Category A). Red and blue horizontal lines at the top of each plot depict orientation channels whose estimated β coefficients are significantly below or above zero, respectively (FDR-corrected permutation test; $p < 0.05$). Shaded regions are ±1 within-participant S.E.M.

**Figure 9. Category Biases Scale Inversely with Distance from the Category Boundary.** (A) The reconstructions shown in Fig. 3 sorted by the absolute angular distance between each exemplar and the category boundary. In our case, the 15 orientations were bisected into two groups of 7, with the remaining orientation serving as the category boundary. Thus, the maximum absolute angular distance between each orientation category and the category boundary was 48°. Participant-level reconstructions were pooled and averaged across visual areas V1, V2, and V3 as no differences were observed across these regions. Shaded regions are ±1 within-participant S.E.M. (B) shows the amount of bias for exemplars located 1, 2, 3, or 4 steps from the category boundary (quantified via a curve-fitting analysis). Error bars are 95% confidence intervals. a.u., arbitrary units.

**Figure 10. Cortical Areas Supporting Robust Decoding of Category Information.** We trained a linear support vector machine to discriminate between activation patterns associated with Category A and Category B exemplars (see *Searchlight Classification Analysis*; Methods). The trained classifier revealed robust category information in multiple visual, parietal, temporal, and prefrontal cortical areas, including many regions previously associated with categorization (e.g., posterior parietal cortex and lateral prefrontal cortex).

**Figure 11. Stimulus Reconstructions in Visual, Parietal, and Frontal cortical areas during the Orientation Mapping and Categorization Tasks.** During the orientation mapping task, participants detected and reported the identity of a target presented in a stream of letters at fixation. During the categorization experiment, participants categorized stimulus orientation into two discrete groups. Shaded regions are ±1 within-participant S.E.M. IPL, inferior parietal lobule; IPS, intraparietal sulcus; sPCS, superior precentral sulcus; IT, inferotemporal cortex, IFG, inferior frontal gyrus. a.u., arbitrary units.

**Figure 12. Summary of Experiment 2.** (A) Participants viewed displays containing an aperture of iso-oriented bars flickering at 30 Hz. (B) The 30 Hz flicker entrained a frequency-specific response known as a steady-state visually-evoked potential (SSVEP). (C) Evoked 30 Hz power was largest over occipitoparietal electrode sites. We computed stimulus reconstructions (Fig. 7) using the 32 scalp electrodes with the highest power. The scale bar indicates the proportion of participants (out of 27) for which each electrode site was ranked in the top 32 of all 128 scalp electrodes. (D-E) Participants categorized stimuli with a high degree of accuracy; incorrect and slow responses were observed only for exemplars adjacent to a category boundary. Shaded regions are ±1 within-participant S.E.M.

1149



1150
1151 **Figure 13. Category Biases Emerge Shortly after Stimulus Onset.** (A) Time-resolved
1152 reconstruction of stimulus orientation. Dashed vertical lines at time 0.0 and 3.0 seconds mark
1153 stimulus on- and offset, respectively. (B) Average channel response function during the first 250
1154 ms of each trial. The reconstructed representation exhibits a robust category bias ($p < 0.01$;
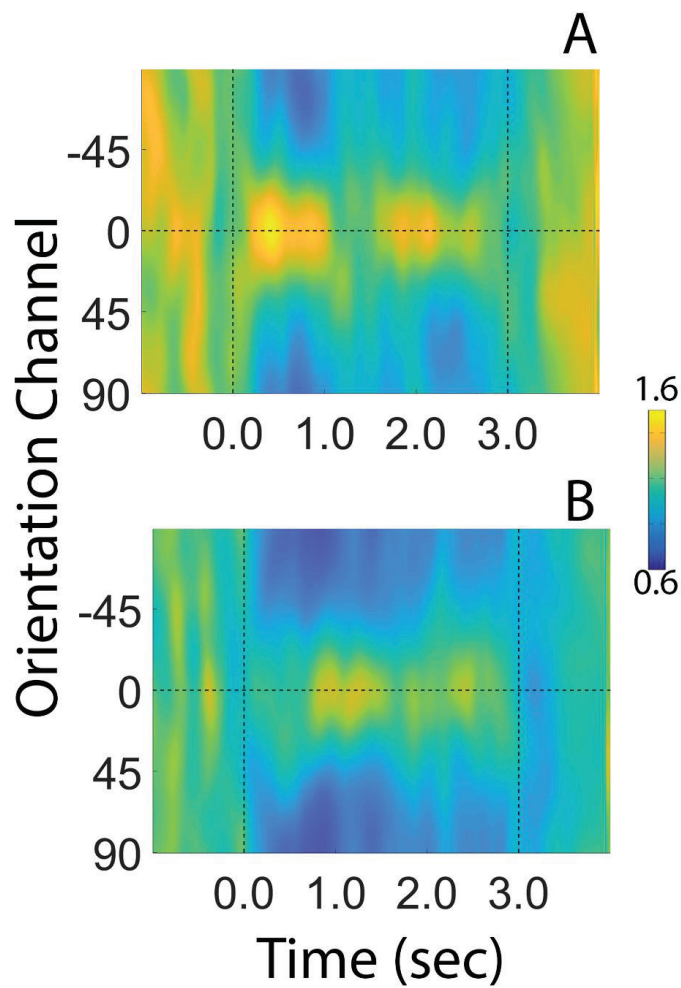1155 bootstrap test). a.u., arbitrary units.
1156
1157

**Figure 14. Stimulus- and category information are absent in pre-trial EEG activity.** Time-averaged reconstruction computed over an interval spanning -250 to 0 ms relative to stimulus onset. The center of the reconstruction was statistically indistinguishable from 0° ($p = 0.234$; bootstrap test)
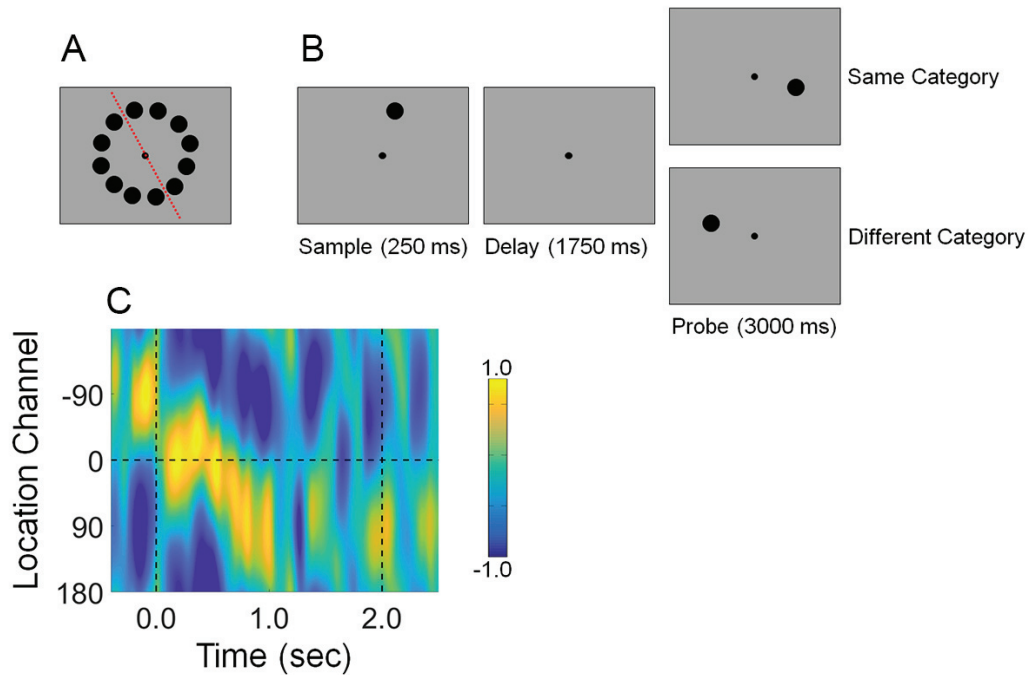
**Figure 15**. **Reconstructions of stimulus orientation during the orientation mapping task (A) and the category discrimination task (B) during Experiment 2**. Vertical dashed lines at time 0.0 and 3.0 mark the start and end of each trial, respectively. Reconstructions were computed using a leave-one-run-out cross validation approach where data from N-1 runs were used to estimate channel weights and data from the remaining run were used to estimate channel responses. This procedure was iterated until all runs had been used to estimate channel responses and the results were averaged over permutations. Units of response are arbitrary.
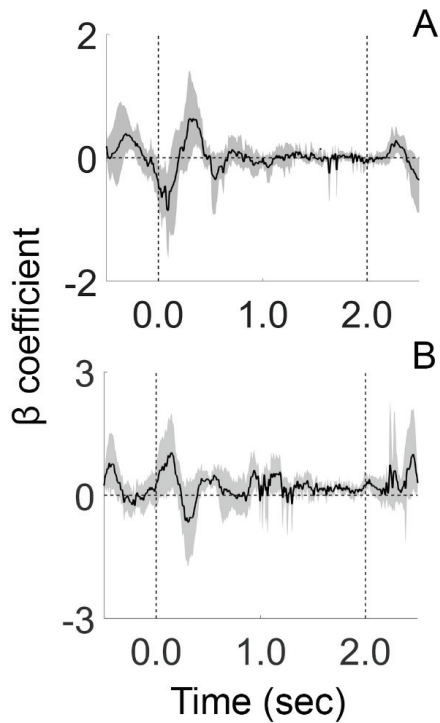
**Figure 16. No systematic biases in eye position during orientation categorization (Experiment 2).** We regressed trial-by-trial records of stimulus orientation (A) or category (B) onto horizontal EOG activity. Thus, positive coefficients reflect a systematic relationship between stimulus orientation (or category) and eye position. No such biases were observed. Black vertical dashed lines at 0.0 and 3.0 depict the start and end of each trial. Shaded regions depict the 95% within-participant confidence interval of the mean.

1205
1206
**Figure 17. Design and Results of Experiment 3.** (A) Possible stimulus locations. The
orientation of the category boundary (red dashed line) was randomly determined for each
participant (example shown). (B) Delayed match-to-category (DMC) task. Participants
remembered the position of a sample disc over a blank delay, then judged whether the location of
a probe disc was drawn from the same location category or a different location category. In this
example, the categories are defined by the boundary shown in panel A. (C) Location-specific
reconstructions computed during the DMC task. Vertical dashed lines at 0.0 and 2.0 sec mark the
onset of the sample and probe epochs, respectively. Participants could not prepare a response
until the onset of the probe display, yet a robust category bias was observed during the delay
period. This suggests that category biases observed in Experiments 1 and 2 are not solely due to
mechanisms of response selection.

1218
1219

1220

**Figure 18. No systematic biases in eye position during location categorization (Experiment 3).** We regressed trial-by-trial records of stimulus location (A) or category (B) onto horizontal EOG activity. Thus, positive coefficients reflect a systematic relationship between stimulus orientation (or category) and eye position. No such biases were observed. Black vertical dashed lines at 0.0 and 3.0 depict the start and end of each trial. Shaded regions depict the 95% within-participant confidence interval of the mean.

1227