1  **Spatial tuning shifts increase the discriminability and fidelity of population codes in visual**
2  **cortex**
3
4  Running title: Spatial attention from units to populations
5

6  Vy A. Vo[1*], Thomas C. Sprague[1,2], John T. Serences[1,3,4*]
7  [1]Neurosciences Graduate Program, University of California, San Diego, La Jolla, CA 92093
8  [2]Department of Psychology and Center for Neural Science, New York University, New York,
9  New York 10003
10  [3]Department of Psychology, University of California, San Diego, La Jolla, CA 92093
11  [4]Kavli Institute for Brain and Mind, University of California, San Diego, La Jolla, CA 92093
12
13  [*]**Correspondence:**
14  Neurosciences Graduate Program
15  University of California, San Diego
16  La Jolla, CA 92093-0634
17  vyaivo@ucsd.edu or jserences@ucsd.edu
18
19

25  **ABSTRACT**

26  Selective visual attention enables organisms to enhance the representation of behaviorally

27  relevant stimuli by altering the encoding properties of single receptive fields (RFs). Yet we know

28  little about how the attentional modulations of single RFs contribute to the encoding of an entire

29  visual scene. Addressing this issue requires (1) measuring a group of RFs that tile a continuous

30  portion of visual space, (2) constructing a population-level measurement of spatial

31  representations based on these RFs, and (3) linking how different types of RF attentional

32  modulations change the population-level representation. To accomplish these aims, we used

33  fMRI to characterize the responses of thousands of voxels in retinotopically organized human

34  cortex. First, we found that the response modulations of voxel RFs (vRFs) depend on the spatial

35  relationship between the RF center and the visual location of the attended target. Second, we

36  used two analyses to assess the spatial encoding quality of a population of voxels. We found that

37  attention increased fine spatial discriminability and representational fidelity near the attended

38  target. Third, we linked these findings by manipulating the observed vRF attentional modulations

39  and recomputing our population measures. Surprisingly, we discovered that attentional

40  enhancements of population-level representations largely depend on position shifts of vRFs,

41  rather than changes in size or gain. Our data suggest that position shifts of single RFs are a

42  principal mechanism by which attention enhances population-level representations in visual

43  cortex.

44

45  **INTRODUCTION**

46  Spatial receptive fields (RFs) are the basis of information processing throughout the visual

47  system. They are directly modified by selective visual attention to improve the fidelity of sensory

48    representations, which likely enables more precise and accurate behavioral choices (Desimone

49    and Duncan, 1995; Anton-Erxleben and Carrasco, 2013). Prior macaque studies have found that

50    covert spatial attention changes the position, size, and amplitude of responses from single-cell

51    RFs in early cortical areas such as V1, V4, and MT by (Moran and Desimone, 1985; Connor et

52    al., 1996, 1997, Womelsdorf et al., 2006, 2008; Roberts et al., 2007; David et al., 2008). Recent

53    neuroimaging studies have also shown that single voxel RFs (vRFs) undergo similar response

54    changes with attention, shifting towards the attended target or changing in size (Klein et al.,

55    2014; Kay et al., 2015; Sheremata and Silver, 2015). Most accounts suggest that these RF

56    changes improve the spatial representations of the attended target, either by boosting the signal-

57    to-noise ratio (SNR) by increasing response amplitude, or by increasing the spatial resolution of

58    the target representation by decreasing the size or tuning width (Desimone and Duncan, 1995;

59    Anton-Erxleben and Carrasco, 2013; Cohen and Maunsell, 2014). These mechanisms are akin to

60    turning up the volume knob (gain increase) or to using smaller pixels to encode a digital image

61    (size decrease).

62        Despite all of these documented modulations, it is not yet clear how these different types

63    of RF modulations are combined to support more robust population codes. Recent studies have

64    only begun to explore how interactions between neurons may affect the coding properties of the

65    population (Anton-Erxleben and Carrasco, 2013; Cohen and Maunsell, 2014). Yet analyzing

66    these data at a population level is crucial for understanding how spatial attention changes the

67    representation of an entire visual scene. Prior fMRI studies that measured many vRFs across

68    space do not report the full pattern of response modulations across space, such as changes in the

69    amplitude or size of vRFs with respect to the attended target (Sprague and Serences, 2013; Klein

70    et al., 2014; Kay et al., 2015). The first aim of this study was to evaluate how all of these

71     properties jointly change across space with spatial attention, in retinotopic regions from V1

72     through posterior IPS.

73          The second aim of the study was to evaluate how different types of RF modulations

74     contribute to population-level enhancements of an attended region of space. Single RFs in early

75     visual areas are fundamentally local encoding models that are relatively uninformative about

76     spatial regions outside their immediate borders. To study their relationship to a population-level

77     representation of space, we also need a larger-scale spatial encoding model which can

78     incorporate information from many underlying, spatially selective RFs to form a representation

79     of the entire visual scene.  Here, we used two different population-level metrics of spatial

80     encoding fidelity to investigate these questions. Specifically, we tested how changes in vRF

81     amplitude, size, or position affect two measurements of population-level representations: (1) the

82     spatial discriminability of population codes and (2) stimulus reconstructions based on a

83     multivariate inverted encoding model.

84          We found that vRF position shifts increase both the spatial discriminability of the

85     attended region as well as the fidelity of stimulus reconstructions near the attended target.

86     Surprisingly, shifts in vRF position captured more of the population-level enhancements with

87     attention than changes in vRF size or gain. This poses problems for traditional 'labeled-line'

88     models of information processing, which posit that each RF in the visual pathway relies on

89     consistently 'labeled' inputs from RFs in an earlier visual area. Our findings suggest that these

90     shifts in the 'labels' of RFs play an important role in the attentional enhancement of visual

91     information, and that labeled-line models may need to be reformulated to accommodate these

92     data.

93

94    **MATERIALS & METHODS**

95    **Task design and participants**

96    We collected data from 9 human participants (4 female), 6 of whom had previously completed a

97    set of retinotopic mapping scans in the lab (participants AA, AB, AC, AI, and AL in Sprague &

98    Serences, 2013; participants AA, AC, and AI in (Sprague et al., 2014); all participants in (Ester

99    et al., 2015). All participants provided written informed consent and were compensated for their

100   time ($20/hour) as approved by the local UC San Diego Institutional Review Board. Participants

101   practiced both the attention task and the localizer task before entering the scanner. A minimum

102   of four hours of scanning was required to complete the entire analysis, so one participant was

103   excluded due to insufficient data (they only completed 2 hours). Another participant was

104   excluded for inconsistent behavioral performance, with average task accuracy at chance (48.6%).

105   This yielded a total of 7 participants who completed the entire experiment (3 2-hour scan

106   sessions per participant).

107          Participants centrally fixated a gray rectangular screen (120x90 cm) viewed via a head-

108   coil mounted mirror (~3.85 m viewing distance). They attended one of three fixed locations on

109   the screen: the fixation point or a target to the left or right. During each 2000 ms trial, subjects

110   reported a change in the attention target. When subjects attended fixation, they reported whether

111   a brief contrast change (100 – 400 ms, starting 300 – 1000 ms into the trial) was dimmer or

112   brighter than the baseline contrast. The peripheral attention targets were two pentagons (0.17°

113   radius; 50% contrast) centered 2.1º to the left and right of fixation (**Fig 1A**). When subjects

114   attended a peripheral target, they reported whether it rotated clockwise or counter-clockwise

115   (rotation duration 100 - 300 ms, starting 300 - 1600 ms into the trial). Inter trial intervals (ITIs)

116   randomly varied between 1000 to 3000 ms in 500 ms increments (mean ITI: 2000 ms). The

117    magnitude of the contrast change or the rotation was adjusted on each run to keep task

118    performance for each participant near 75% (mean = 75.90%, bootstrapped 95% C.I. [72.46%,

119    79.20%]), with no significant difference between conditions as evaluated with a one-way

120    repeated measures ANOVA randomization test ($F(1,11) = 0.220$, randomized $p = 0.802$). For

121    four participants, we collected 6 runs on the attend periphery tasks without a change in the

122    luminance of the fixation stimulus. Performance on the attend periphery tasks was stable across

123    runs with and without the luminance change (repeated-measures ANOVA with run type x

124    random participants factor; $p = 0.439$, null F distribution using randomized labels for 10,000

125    iterations). Therefore, these data were collapsed across scan sessions with and without changes

126    in fixation luminance.

127          On 51 of the 61 trials in each run, a full-contrast 6 Hz flickering checkerboard (0.68°

128    radius; 1.67 cycles/deg) appeared for 2000 ms at one of 51 different locations across the screen

129    to map the spatial sensitivity of visually responsive voxels. When one of these checkerboards

130    overlapped with any of the static attention targets, they were partially masked with a small

131    circular aperture the same color as the screen background (0.16°/0.25° radius aperture for

132    fixation/pentagon, respectively) that allowed the stimulus to remain seen. Participants were

133    instructed to ignore the task-irrelevant flickering checkerboards throughout the experiment.

134    During the 10 null trials, the participant continued to perform the attention task but no

135    checkerboard was presented. Null trials and mapping stimulus trials were presented in a

136    psuedorandom interleaved order.

137          The location of the checkerboard mapping stimulus on each trial was determined by

138    generating an evenly spaced triangular grid (0.84° between grid points) and centering the

139    checkerboard on one of these grid points. The location of the checkerboard mapping stimulus

140  was then jittered a random amount from these grid points (+/- 0.42°/0.37° horizontal/vertical).

141  When subjects attended the peripheral target, half of the runs were presented at the discrete grid

142  positions so that we could achieve more stable stimulus reconstructions (see **Population analysis**

143  **(2)**).

144  **Magnetic resonance imaging**

145  We obtained all structural and functional MR images from participants using a GE 3T MR750

146  scanner at University of California, San Diego. We collected all functional images (19.2 cm$^2$

147  FOV, 64 x 64 acquisition matrix, 35 interleaved slices, 3 mm$^3$ voxels with 0 mm slice gap, 128

148  volumes per scan run) using a gradient echo planar pulse sequence (2000 ms TR, 30 ms TE, 90°

149  flip angle) and a 32-channel head coil (Nova Medical, Wilmington, MA). Five dummy scans

150  preceded each functional run. A high-resolution structural image was acquired at the end of each

151  session using a FSPGR T1-weighted pulse sequence (25.6 cm$^2$ FOV, 256 x 192 acquisition

152  matrix, 8.136/3.172 ms TR/TE, 172 slices, 9° flip angle, 1 mm$^3$ voxels, 192 volumes). All

153  functional scans were co-registered to the anatomical images acquired during the same session,

154  and this anatomical was in turn co-registered to the anatomical acquired during the retinotopy

155  scan.

156      EPI images were unwarped with a custom script from UCSD's Center for Functional

157  Magnetic Resonance Imaging using FSL and AFNI. All subsequent preprocessing was

158  performed in BrainVoyager 2.6.1, including slice-time correction, affine motion correction, and

159  temporal high-pass filtering to remove slow signal drifts over the course of each session. Data

160  were then transformed into Talairach space and resampled to 3x3x3 mm voxel size. Finally, the

161  BOLD signal in each voxel was transformed into Z-scores on a scan-by-scan basis. All

162    subsequent analyses were performed in MATLAB using custom scripts (to be available online

163    upon acceptance).

164         We constrained our analyses to visually responsive voxels in occipital and parietal cortex

165    using a separate localizer task (3-5 runs per participant). On 14 trials, participants fixated

166    centrally and viewed a full-field flickering checkerboard (10 Hz, 11.0/8.3° width/height) for

167    8000 ms. Participants detected whether a small area (2D Gaussian, $\sigma = 0.2°$) within the

168    checkerboard dimmed in contrast. Contrast dimming occurred between 500 to 4000 ms after the

169    start of the trial, and lasted between 2000 to 3000 ms (all uniformly sampled in 500 ms steps).

170    This contrast change occurred infrequently (randomly on 5 out of 14 trials) at a random location

171    within the checkerboard. The average contrast change was varied between runs to maintain

172    consistent performance at ~75% accuracy (mean performance 78.0%). On 8 trials participants

173    simply fixated throughout the trial without a checkerboard being presented. For all subsequent

174    analyses, only voxels in the retinotopically defined areas V1, V2, V3, V4, V3A/B and IPS0 with

175    a significantly positive BOLD response to the localizer task (at FDR $q = 0.05$) were included

176    (Benjamini and Yekutieli, 2001).

177         For all subsequent analyses, we used trial-wise BOLD z-scores. We estimated these by

178    creating an event predictor for each checkerboard mapping stimulus and convolving it with a

179    canonical two-gamma HRF (peak at 5 s, undershoot peak at 15 s, response undershoot ratio 6,

180    response dispersion 1, undershoot dispersion 1). We then solved a general linear model (GLM)

181    to find the response to each predictor. To standardize our data across runs, we z-scored the

182    BOLD responses within each run and concatenated the z-scores across runs.

183    **Statistical procedures**

184     All reported confidence intervals (CIs) are computed by resampling the data with replacement

185     (i.e. bootstrapping). The number of iterations for each bootstrapping procedure varied

186     (depending on available computing power and time for that procedure) and are therefore reported

187     with each result. For tests comparing a bootstrapped distribution against zero, p-values were

188     computed by conducting two one-tailed tests against 0 (e.g., mean(param_change < 0) &

189     mean(param_change > 0)) and doubling the smaller p-value. All repeated tests were FDR

190     corrected (q = 0.05).

191

192     **Voxel receptive field (vRF) estimation, fitting, and parameter analysis**

193         We first estimated vRFs for each attention condition to investigate (1) how vRF

194     parameters changed when participants attended to different locations and (2) the spatial pattern

195     of vRF changes across visual space. We note here that prior reports have referred to similar

196     voxel RF models as population receptive fields, or pRFs, to emphasize the fact that each voxel

197     contains a population of spatially tuned neurons (Dumoulin and Wandell, 2008; Wandell and

198     Winawer, 2015). However, since we are comparing modulations at different scales in the present

199     study (i.e. modulations in single voxels and in patterns of responses across many voxels), we will

200     refer to these single voxel measurements as voxel receptive fields (vRFs), and will reserve the

201     term 'population' exclusively for multivariate measures involving several voxels, allowing our

202     terminology to align with theories of population coding.

203         We estimated voxel receptive fields (vRFs) using a modified version of a previously

204     described technique (Sprague and Serences, 2013). This method estimates a single voxel's

205     spatial sensitivity by modeling its BOLD responses as a linear combination of discrete, smooth

206     spatial filters tiled evenly across the mapped portion of the visual field. These spatial filters (or

207    spatial channels) form our modeled basis set. We then regressed the BOLD z-scores ($v$ voxels x

208    $n$ trials) onto a design matrix with predicted channel responses for each trial ($C$, $k$ channels x $n$

209    trials) by solving Equation 1:

210    (1)      $B = WC$

211    for the matrix $W$ ($v$ voxels x $k$ channels).

212          Each of the $k$ channels in the basis set was defined as a two-dimensional cosine that was

213    fixed to reach 0 at a set distance from the filter center:

214    (2)      $f(r) = \left( 0.5 \left( \cos\left(\frac{r\pi}{s}\right) + 0.5 \right) \right)^7 \; for \; r < s,$

215    where $r$ is the distance from the filter center and $s$ is the size constant. Setting a zero baseline in

216    this function ensured that we could estimate a stable baseline for each voxel by restricting the

217    response of the channel to a known subregion of the visual display. Since the estimated vRF size

218    depends on the size of the filters, we made the filters fairly small (1.08° FWHM) and dense (91

219    filters arranged in a 13 horizontal / 7 vertical grid, each spaced 0.83° apart). We then discretized

220    the filters by sampling them in a high-resolution 2D grid of 135 by 101 pixels spanning 10° by

221    5°. The discretized filters ($k$ filters by $p$ pixels) were multiplied with a mask of the checkerboard

222    stimulus on every trial ($p$ pixels by $n$ trials) so that the design matrix $C$ contained predictions of

223    the spatial channel responses on every trial of the mapping task.

224          In order to fit our estimated vRFs with a unimodal function, we used ridge regression to

225    solve Equation 1. This is a common regularization method which sparsifies the regression

226    solution by penalizing the regressors with many small weights (Lee et al., 2013). This meant

227    solving for an estimate of $W$ by the following:

228    (3)      $\widehat{W}^T = (CC^T + \lambda I)^{-1}CB^T,$

229    where λ is the ridge parameter penalty term, and $I$ is a $k$ x $k$ identity matrix. We estimated an

230    optimal λ for each voxel by evaluating Equation 3 over a range of λ values (0 to 750) for all runs

231    of the attention task (e.g., concatenating all attention conditions together). We then computed the

232    Bayesian Information Criterion (BIC) for each of these λ values, estimating the degrees of

233    freedom in the ridge regression as $df = trace(C \ (C^T C + \lambda I)^{-1} C^T)$. The λ with the smallest BIC

234    was selected for each voxel. Since the attention comparisons are done within voxels, the varying

235    λ penalty across voxels could not explain the attention effects we observed.

236          To select reliable voxels for analysis, we next implemented a set of conservative

237    thresholding steps (**Table 1**). We first needed to select voxels with reliable visual responses, so

238    we only kept voxels with trial beta weights that predicted at least 50% of the BOLD time courses

239    in each scan session. Second, we only used voxels that could be successfully regularized with

240    ridge regression. Any voxels with the maximum λ (750) were discarded, as this indicated that the

241    ridge regression solution had not converged. Finally, we verified that the resulting regression

242    model could predict an independent dataset, so we performed exhaustive leave-one-run-out cross

243    validation for each attention condition. This ensured that the λ estimated across attention

244    conditions produced reliable data for each condition separately. We estimated $W$ using data from

245    all but one run (Equation 3) and used this to predict the BOLD GLM trial estimate of the left-out

246    run (Equation 2), again all done separately for each condition. We then computed the mean

247    correlation between the predicted & real BOLD GLM trial estimates across cross-validation

248    folds for each voxel. Note that it is not possible to calculate a coefficient of determination on

249    regularized data, since the process of ridge regression changes the scale of the predicted data (see

250    Huth et al., 2012 for more). We only kept voxels where this cross-validation r > 0.25 for all 3

251    conditions.

252        To quantify each vRF, we fit each voxel with a smooth 2D function with 4 parameters:

253    center, size, baseline, and amplitude (**Fig 1b**; Equation 2). Here, we define the vRF baseline as

254    the voxel's typical response that does not depend on the position of the mapping stimulus. The

255    vRF amplitude is defined as the spatially-selective increase in a voxel's response above this

256    baseline. Together, these two parameters index how much of the voxel's response is due to a

257    change in mapping stimulus position. Finally, the size and location parameters estimate the

258    spatial selectivity and the spatial position preference of the vRFs, respectively. We first

259    downsampled the vRFs by multiplying the estimated weights $\widehat{W}$ for each voxel (a 1 x $k$ channel

260    vector) by a smaller version of the spatial grid that contained the basis set (68 by 51 pixel grid;

261    10° by 5°). This speeded up the process of fitting the pixelwise surface with Eq. 2. This fitting

262    process began with a coarse grid search that first found the best fit in a discrete grid of possible

263    vRF parameters (center sampled in 1° steps over the mapped portion of the visual field; size

264    constant logarithmically sampled at 20 points between 2.3 and 38.5, which gives FWHMs

265    between 0.9° and 15.3°). We then estimated the best fit amplitude and baseline for each of the

266    grid points using linear regression. The grid point fit with the smallest root mean square error

267    (RMSE) provided the initialization seed to a continuous error function optimization algorithm

268    (fmincon in MATLAB). This fit had several constraints: the final solution must place the center

269    within 2 grid points of the seeded fit (parameterized by position and size) and within the mapped

270    visual field; the amplitude must be between 0 and 5; the baseline must be between -5 and 5

271    BOLD z-score units. Occasionally, this nonlinear fitting algorithm did not converge and resulted

272    in a larger error than the original seed. In this case we took the best fit grid point as the final fit.

273        To test whether vRF fit parameters changed when participants focused spatial attention at

274    different positions, we compared fits during each attend periphery condition with fits during the

275    attend fixation condition. We computed a difference score (attend peripheral – attend fixation) to

276    describe the magnitude of the attentional modulation. For example, a difference score of –2° in

277    the FWHM of the vRF would indicate that the response profile width decreased when the

278    participant attended to the peripheral target. We then tested whether the vRF parameter

279    difference scores differed significantly from 0 within a visual region of interest (ROI) by

280    bootstrapping the distribution of difference scores across participants 10,000 times.

281         To test whether these vRF changes were modulated by their position in the visual field,

282    we first calculated each vRF's distance from the attended location (*v_dist_attn*) using its position

283    during the fixation task. These were sorted into distance bins (0° to 2.5°, in 0.25° steps) and all

284    vRF difference scores in that bin were averaged across participants. We then fit an *n*th order

285    polynomial to the binned vRF difference scores as a function of *v_dist_attn*, where $n = 0, 1,$ or $2$.

286    This corresponds to a constant offset ($0^{th}$ order), a linear fit ($1^{st}$ order), or a quadratic or parabolic

287    fit ($2^{nd}$ order). The most parsimonious fit was chosen using a nested F-test. Fit coefficient CIs

288    were generated by bootstrapping the data across participants 5,000 times before repeating the

289    binning, averaging, and fitting procedure.

290

291    **Population analysis (1): Fine spatial discriminability metric**

292    To compute the spatial discriminability of a population of vRFs, we estimated the spatial

293    derivative of each vRF at every point in the mapped visual field in 0.1° steps (**Fig 1C**). This was

294    done by taking the slope of the vRF along the x and y direction at each pixel in the image of the

295    visual field and squaring this value (Scolari and Serences, 2009, 2010). This measurement is a

296    good descriptor of how well a population code can discriminate small changes in the spatial

297    arrangement of the stimulus array, which depends on the rising and falling edges of a tuning

298     curve rather than the difference between the peak response and a baseline response (Regan and

299     Beverley, 1985; Pouget et al., 2003; Butts and Goldman, 2006; Navalpakkam and Itti, 2007;

300     Scolari and Serences, 2009, 2010). In order to restrict our measurements to the relevant area near

301     the peripheral target, we computed discriminability values within 1 degree of the center of each

302     target across both spatial dimensions (x and y). These were summed and divided by the

303     maximum discriminability value in that population in order to make the results comparable

304     despite changes in vRF coverage or responsiveness.

305

306     **Population measurements (2): Stimulus reconstructions using an inverted spatial encoding**

307     **model**

308     In addition to computing the discriminability metric described above, we also reconstructed an

309     image of the entire visual field on each trial using a population-level encoding model. Compared

310     to the local spatial discriminability index, this is a more sensitive method of assessing the amount

311     of spatial information encoded in a population of voxels because it exploits the pattern of

312     response differences across voxels, rather than treating each voxel as an independent encoding

313     unit (Serences and Saproo, 2012; Sprague et al., 2015).

314             We train the spatial encoding model using a procedure similar to the vRF estimation

315     analysis described above (**Fig 4a**). This yields an estimated matrix of weights, $\widehat{W_2}$ , which

316     specifies how much each voxel in a region of interest responds to each of the spatial channels

317     (Brouwer and Heeger, 2009; Serences and Saproo, 2012; Sprague and Serences, 2013; Sprague

318     et al., 2015). We then solve Eq. 1 using the Moore-Penrose pseudoinverse with no regularization:

319             (4)     $\widehat{W_2} = BC^T(CC^T)^{-1}$

320    *C* was constructed using a set of 48 evenly tiled spatial filters (Eq. 2; 8 horizontal / 6 vertical;

321    spaced 1.43° apart; 1.78° FWHM). $\widehat{W_2}$ was estimated using the data from the jittered position

322    runs. This was done separately for each participant, using a balanced training set that contained

323    an equal number of attend left, attend right, and attend fixation runs.

324          To reconstruct a representation of the mapped visual space, we inverted the model and

325    multiplied the pseudoinverse of the estimated weight matrix $\widehat{W_2}$ with a test dataset from the

326    discrete position runs ($B_2$), yielding estimated channel activations for each trial ($C_2$; $k_2$ channels

327    by $t$ test trials) (Equation 5). Thus, we refer to this analysis as the inverted encoding model

328    (IEM).

329    (5)      $\hat{C}_2 = \left( \widehat{W_2}^T \widehat{W_2} \right)^{-1} \widehat{W_2}^T B_2$

330    Because of mathematical constraints on inverting $W_2$ (number of voxels must be greater than

331    number of channels), we included all voxels in each ROI instead of just the subset of well-fit

332    voxels used in the vRF analyses described above. We performed this inverting procedure twice

333    using different test datasets, once for the discrete position attend left runs and once for the

334    discrete position attend right runs.

335          When we multiply the resulting channel activations by a grid of pixels that define the

336    spatial channels, we obtain a spatial representation of the entire visual field on each trial. This

337    image contains a stimulus reconstruction showing where the checkerboard should have been

338    given the trained model and the activation pattern across all voxels in the independent test set.

339    The stimulus reconstructions were then fit in the same manner as the vRFs, using Eq. 1 to

340    estimate the center, size, amplitude, and baseline (mean fit RMSE across all ROI reconstructions

341    0.173; 95% CI [0.102, 0.312]). Here, the baseline is an estimate of the multivariate

342    reconstruction that is spatially non-selective—i.e., not significantly modulated by the position of

343    the mapping stimulus. The amplitude describes the maximal increase in that reconstruction

344    relative to baseline when the mapping stimulus is on the screen.

345         To assess how attention changed reconstructions of the mapping stimulus across the

346    visual field, we first computed a difference score that described the effect of attention by folding

347    the visual field in half (i.e. collapsing across hemifield) and comparing parameters in the

348    attended vs. ignored hemifield. We excluded the reconstructions that fell along the vertical

349    meridian (3 of 51 stimulus positions). This allowed us to control for the overall effect of

350    eccentricity while remaining sensitive to other spatial patterns in stimulus reconstruction

351    modulations.

352         We then set up a single factor repeated measures omnibus ANOVA to determine which

353    pairs of ROI and parameter (e.g., V1 size, V1 amplitude, etc.) were affected by either attention or

354    Euclidean distance from the target stimuli. The attention factor had two levels (attend/ignore)

355    and the distance factor had 6 levels (6 evenly spaced distance bins from 0° to 2.54°). Based on

356    the results of this omnibus test, we tested any significant ROI-parameter combination in a 2-way

357    repeated measures ANOVA of attention by distance. To estimate the p-values for these tests, we

358    generated empirical null distributions of the F-scores by randomizing the labels within each

359    factor 10,000 times within each participant. Reported p-values are the percentage of the

360    randomized F-scores that are greater than or equal to the real F-scores.

361

362    **Population analysis (3): Layered spatial encoding model to link vRFs to multivariate**

363    **stimulus reconstructions**

364    In order to test how changes in the response properties of the underlying vRFs contributed to

365    changes in the fidelity of region-level stimulus reconstructions, we generated simulated patterns

366    of voxel activity on every trial by predicting the response to each stimulus based on the vRF fit

367    parameters. We then used this simulated data to estimate and invert a population-level spatial

368    encoding model, as described above (**Fig 6A**). Note that for these simulations, we could only use

369    well-fit voxels to generate simulated BOLD timeseries. Therefore, we could not accurately

370    estimate reconstructions for some participant-ROI pairs that had an insufficient number of

371    voxels. Pairs that indicated a poorly conditioned matrix (e.g., number of voxels is fewer than the

372    number of channels) were excluded (total 10 out of 35 pairs; V3 (AI, AP); V3A/B (AL, AP, AT);

373    V4 (AA, AT); IPS0 (AA, AR, AU)).

374        To simulate each voxel's BOLD response on every trial that the participant completed in

375    the real experiment, we first created a high-resolution set of spatial channels (21 by 11 channels

376    spaced $0.5°$ apart, FWHM $= 0.65°$) and generated weights for each channel based on the vRF fit

377    obtained from prior analysis. That is, we evaluated Eq. 2 for each channel at the vRF's fit center

378    and adjusted the response gain by multiplying this result by the fit amplitude and adding the fit

379    baseline. We then added independent Gaussian noise to each of these channel weights,

380    simulating a small amount of variance in the voxel's response ($\sigma = 0.5$). Each voxel's channel

381    weights were then multiplied by the idealized channel response on each trial (that is, the channel

382    filter convolved with the stimulus mask), effectively simulating the BOLD response on each trial

383    for the entire population of voxels according to their measured vRFs. We added Gaussian noise

384    to this simulated response as well ($\sigma = 0.5$). We then computed stimulus reconstructions using

385    the same method as described above (the IEM in **Population analysis (2)**), averaging resulting

386    reconstructions across participants and like positions before fitting.

387        To assess the stability of the reconstructions that were based on simulated data, we

388    repeated the simulations 100 times and averaged across the fits of all iterations to generate the

389  plots in **Fig 6**. Note also that the chosen level of noise did not qualitatively impact the results.

390  For example, rather than just adding Gaussian noise, we also created a noise model that followed

391  the covariance structure between all voxels in each ROI. To estimate the covariance matrix, we

392  computed the residuals between the true trial-wise beta weights and the predicted trial-wise beta

393  weights for each voxel based on its vRF model. We then calculated the pairwise covariance

394  between the residuals for each set of voxels. Last, we added noise that followed this covariance

395  structure to each voxel's channel weights and simulated BOLD response. This noise was scaled

396  to be the same as the noise level that most accurately captured the real reconstruction data (i.e.,

397  mean noise is 0.5 standard units). The pattern of results between each of the model

398  manipulations remained the same, so those results are not discussed here.

399      To compare whether the results of the layered model differed significantly from the

400  reconstructions generated with real data, we first calculated difference scores across attention

401  condition (attended – ignored; see **Population analysis (2)**). This yielded 24 difference scores

402  each for both attention conditions. Since the real data did not have any repeated iterations, we

403  averaged across all 100 iterations of the model to match the dimensionality of the real

404  reconstructions (2 conditions x 24 difference scores x 4 parameters). We then calculated the error

405  between the difference scores from the full empirical dataset (as the data shown in **Fig 5**) and the

406  modeled data to obtain the root mean square error (RMSE).

407      To test how shifts in vRF centers contributed to changes in the stimulus reconstructions,

408  we also generated reconstructions from modeled voxels that had the same fit center across both

409  attention conditions. We defined each voxel's vRF center as the fit center obtained from that

410  voxel during the neutral attend fixation condition. We then repeated the stimulus reconstruction,

411  reconstruction fitting, and statistical testing as described above. A similar procedure was

412    repeated for all reported combinations of parameter changes across conditions. Again, whichever

413    parameter was held constant took its value from the neutral attend fixation condition.

414        To calculate the confidence intervals on the RMSE changes in **Fig 6C**, we resampled

415    with replacement across the 100 model iterations and took the difference between the RMSE of

416    the null model, in which no parameters varied between attention conditions, and the RMSE of

417    the model which held some number of vRF parameters constant across attention conditions. This

418    procedure was repeated 500 times.

419

420    **RESULTS**

421    **Modulations of vRF properties with spatial attention**

422    We estimated single voxel receptive fields (vRFs) for each voxel in 6 retinotopically-identified

423    visual areas from V1 to IPS0. The estimation of vRFs was done independently for each attention

424    condition so that we could compare a single voxel's spatial tuning across conditions.

425        To confirm that the fit sizes were consistent with previous results, we fit a line to the

426    estimated sizes as a function of the vRF center eccentricity. First, we combined all vRFs across

427    participants and conditions in each ROI. We then binned the vRF centers every 0.25° from

428    fixation and calculated the mean size (**Fig 2b**). We first replicated an increase in vRF size with

429    increasing eccentricity, and an increase in the slope of this relationship across visual regions

430    (Gattass et al., 2005; Dumoulin and Wandell, 2008; Amano et al., 2009; Harvey and Dumoulin,

431    2011) (**Fig 2b**). These observations confirm that our method produced reasonable vRF estimates

432    that were consistent with previous reports.

433        Covert attention to either the left or the right position modulated vRF properties by

434    shifting vRFs significantly closer to the attended location compared to the attend fixation

435    condition ($p < 0.005$ in all ROIs). While we did observe size changes in individual voxels, the

436    mean change was not significantly different from zero ($p > 0.05$ in all ROIs). Size increases have

437    been previously reported in tasks that required subjects to attend to the mapping stimulus rather

438    than to ignore it, as in the present study (Sprague and Serences, 2013; Kay et al., 2015;

439    Sheremata and Silver, 2015). In these previous studies, the locus of attention changes on each

440    trial. Accordingly, if attention attracts RFs as our data suggest, these combined shifts might

441    average out to form a larger RF estimate. Here, we fixed the locus of attention so we could more

442    finely characterize the effects of focal attention and found no net change in vRF size. However,

443    we did find a general increase in vRF response gain, such that amplitude increased while mean

444    baseline decreased ($p < 0.001$ for all tests). Since all of these measures were calculated relative to

445    a fixation task, these data suggest that covert spatial attention to a peripheral location caused

446    widespread position and gain modulations in all vRFs across the visual field.

447            Previous reports in humans and monkeys have suggested that the preferred position of

448    RFs shift when subjects covertly attend to an area in the visual field (Womelsdorf et al., 2006,

449    2008; Klein et al., 2014) and when they make visually-guided saccades to an attended location

450    (Zirnsak et al., 2014). It is unclear, however, whether these position shifts all radially converge

451    towards the attended target or whether the RFs shift uniformly along a vector extending from

452    fixation to the attention or saccade target (Tolias et al., 2001; Klein et al., 2014; Zirnsak et al.,

453    2014). Furthermore, reports of other RF properties (such as size) modulating with attention have

454    been mixed (Connor et al., 1996, 1997; Womelsdorf et al., 2008; Niebergall et al., 2011; Sprague

455    and Serences, 2013; Klein et al., 2014; Kay et al., 2015; Sheremata and Silver, 2015). We

456    therefore examined whether each of the vRF parameter changes was dependent on the vRF's

457    location in the visual field, relative to the attended location. First, we created radial distance bins

458     centered on the left or right attended locations, and sorted voxels into these bins based on their

459     preferred position during the fixation condition. After this sorting procedure, data from the right

460     condition were flipped and collapsed with the left condition.

461         When we plotted vRF position changes in each bin, we found that spatial attention caused

462     vRF position shifts that converged on the attended location (two-tailed sign test on vector

463     direction, p < .001 in all ROIs). That is, vRFs shifted closer to the attended location (**Fig 2c**),

464     compared to when subjects attended fixation (mean shift across all vRFs and ROIs: -0.244,

465     bootstrapped 95% C.I. [-0.607, -0.038], **Fig 2d**). Note that small eye movements toward the

466     attended location cannot explain receptive field convergence: this would cause all vRFs to shift

467     in the same horizontal direction, rather than radially converging on one point. These data are

468     consistent with results from both humans (Klein et al., 2014) and macaques (Connor et al., 1996,

469     1997, Womelsdorf et al., 2006, 2008). However, the prior study in humans focused only on vRFs

470     with preferred locations that were foveal to the attended location, and the studies in macaques

471     only report RF position changes in V4 and MT. By contrast, we show that vRF centers converge

472     on the attended location across all visual areas, including primary visual cortex, as well as in

473     vRFs with centers peripheral to the attended location.

474         These plots (**Fig 2a, 2d**) also suggested that vRFs farther from the attended location

475     underwent larger position changes with covert shifts of spatial attention. The size of the

476     attentional modulation may be dependent on the distance between the vRF center and the

477     attended target. To test for this, we fit a polynomial to the vRF parameter changes as a function

478     of distance from the attended location (**Materials and Methods**). We selected the most

479     parsimonious fit ranging from a mean change in vRF parameter ($0^{th}$ order polynomial) to a

480    parabolic change (2nd order polynomial) by conducting a nested F-test (**Table 2**). The significant

481    polynomials of order n > 0 are plotted in **Fig 2e**.

482         Position changes with attention were significantly modulated by the initial vRF position

483    (relative to the locus of attention) in higher visual areas such as V4, as indicated by a significant

484    slope coefficient from a linear fit. While all size changes were best described by a quadratic

485    function, only area V2 showed a significant modulation: we observed size increases in V2 vRFs

486    near the attended location, size decreases about 1 degree away from the attended location, and

487    size increases again for vRFs farther away. Amplitude changes with attention were significant in

488    visual areas higher up in the visual hierarchy, namely V3A/B and IPS0. In these visual areas,

489    amplitude increased for vRFs farther from the attended location. This was paired with a large

490    mean decrease in baseline (**Fig 2d**). These tests suggest that the spatial relationship between the

491    vRF and the attended target changes the type and magnitude of the attentional modulation in

492    different visual areas, consistent with findings from macaque neurophysiology (Connor et al.,

493    1996; Niebergall et al., 2011).

494    **Increases in spatial discriminability depend primarily on vRF position shifts**

495         Next, we assessed how different types of RF modulations influenced the precision of

496    population-level codes for spatial position. We first computed a discriminability metric that

497    described the ability of a population of tuning curves to make fine spatial judgments (**Materials**

498    **and Methods**). Spatial discriminability near the attended target increased relative to the ignored

499    target in the opposite visual hemifield (**Fig 3a**).

500         We then tested how different types of vRF modulations (such as size changes or position

501    shifts) affected this spatial discriminability metric. As a baseline comparison, we first computed

502    discriminability based on vRFs estimated during the attend fixation runs for each subject. We

503    then added different sets of observed attentional modulations to the population before

504    recomputing spatial discriminability. For example, we shifted all the vRF centers to match the

505    measurements when a subject was attending to the left target and computed discriminability near

506    the attended target. Since the response baseline of a vRF does not affect the discriminability

507    metric, we excluded this type of attentional modulation from these analyses.

508          Across all ROIs, we found that vRF position shifts played the biggest role in increasing

509    fine spatial discriminability compared to changes in size or changes in amplitude (**Fig 3b**).

510    Position modulations alone led to a large increase in spatial discriminability, while other

511    combinations of parameter modulations only had an impact if we added in position shifts (i.e. a

512    change in size and position increased discriminability, but size alone did not). The only departure

513    from these patterns was observed in IPS0, where a combination of amplitude and size produced

514    an increase in discriminability even in the absence of changes in vRF position.

515    **Spatial attention increases the fidelity of population-level stimulus reconstructions**

516          By design, the spatial discriminability metric we computed is only informative about

517    local spatial representations, and cannot assess how different patterns of vRF modulations might

518    result in representational changes across the visual field. To address this point, we built a

519    multivariate spatial encoding model to measure how attention changes the representations of

520    visual information in disparate parts of space. This also allowed us to further test the effects of

521    vRF modulations on the encoding properties of the population, including response baseline

522    changes that were not captured by our discriminability metric.

523          The spatial inverted encoding model (IEM) reconstructed an image of the entire visual

524    field on each test trial. We first trained the model using the responses of each voxel on a set of

525    training trials with known mapping stimulus positions. We then created image reconstructions on

526  independent test trials by inverting the model and multiplying it by the voxel responses during

527  each test trial (**Fig 4a**; **Materials and Methods**). Each image contained a representation of

528  where the mapping stimulus should have been given the pattern of voxel activations on that

529  particular trial. The IEM successfully reconstructed the task-irrelevant mapping stimuli using

530  activation patterns across voxels in each visual area from V1 through IPS0 (**Fig 4b**; grand mean

531  error between fit and actual position 2.36°, 95% CI [0.56°, 4.89]).

532       We used these stimulus reconstructions as a proxy for the quality of the spatial

533  representations encoded in a population of voxels. This is line with previous studies showing that

534  stimulus reconstructions have change in amplitude or size as a function of cognitive demands.

535  (Brouwer and Heeger, 2013; Ester et al., 2013; Sprague and Serences, 2013; Sprague et al., 2014,

536  2015, 2016).

537       First, we compared how reconstructed representations of each mapping stimulus changed

538  as subjects shifted their spatial attention. We ran a repeated measures ANOVA of *attention* x

539  *distance bin* for each reconstruction fit parameter (**Materials and Methods**). Here, a main effect

540  of attention would suggest that stimulus reconstructions in the attended hemifield changed in a

541  consistent way compared to the ignored hemifield. A main effect of distance would suggest that

542  stimulus reconstruction changes had a consistent spatial pattern across both the attended and

543  ignored hemifields. This would occur when a stimulus' representation was altered with distance

544  from the attention target. For example, the stimulus reconstruction center should vary linearly

545  with the stimulus' true distance from the attention target. And lastly, an interaction effect would

546  suggest that the distance effect was dependent on whether the reconstruction belonged to the

547  attended or ignored hemifield. In our task, the reconstructed stimuli are always irrelevant to the

548  task of the observer. We therefore predicted an interaction effect where spatial attention would

549    selectively modulate stimulus reconstructions around the attended location (Connor et al., 1996,

550    1997).

551            We found that reconstruction amplitude was selectively increased near the attended

552    location in V3, V4, V3A/B, and IPS0 (interaction effect, bootstrapped p < 0.005; **Fig 5**; **Table**

553    **3**). This can be interpreted as a local boost in SNR. Prior reports found that attending to the

554    mapping stimulus – as opposed to attending to a peripheral target as in the current experiment –

555    caused an increase in the amplitude of all stimulus reconstructions (Sprague and Serences, 2013).

556    That is, representations of task-relevant stimuli increased in SNR. We find here that even

557    representations of task-*irrelevant* stimuli near the attended region of space increase in amplitude,

558    consistent with the idea of an attentional 'spotlight' which boosts the fidelity of spatial

559    representations near the attention target.

560            Although the amplitude interaction effect was present in most visual areas we tested (**Fig**

561    **5**), we found other effects limited to V3A/B and IPS0 that involved modulations in stimulus

562    representations in the ignored hemifield. In these regions, we found that stimulus reconstructions

563    in the ignored hemifield shifted away from the ignored target location (interaction, bootstrapped

564    p < 0.005). We also observed a relative size increase near the ignored attention stimulus in IPS0

565    (interaction, bootstrapped p = 0.005). These results suggest that stimulus reconstructions in the

566    ignored hemifield are less spatially precise in posterior parietal cortex. Finally, there was also a

567    main effect of attention on reconstruction size and baseline in areas V4 & V3A/B (bootstrapped

568    p's <= 0.005). However, unlike the interaction effect in IPS0, size changes in V4 and V3A/B did

569    not vary as a function of distance between the reconstruction and the attended target location.

570    **Using a layered encoding model to explore how single voxel RFs change population-level**

571    **codes**

572    In our final analysis, we used a layered spatial encoding model to determine how changes

573    in vRF properties affected the representations of mapping stimuli in the multivariate

574    reconstructions discussed in the previous section (**Fig 1c**; **Fig 4a**). The goal of this analysis was

575    to determine which vRF modulations contribute the most to the observed increase in the

576    amplitude of stimulus representations around the attended location (**Fig 5**). This analysis thus

577    complements our analysis of the spatial discriminability metric which demonstrated that vRF

578    position changes significantly increased the ability of the population to make fine spatial

579    discriminations near the attention target (**Fig 3c**).

580    The layered spatial encoding model we built links the response properties of single

581    voxels to the encoding properties of a whole population of voxels in a region of visual cortex

582    (**Fig 6a**). In the first layer of the model, we used the fit vRFs to generate simulated BOLD data

583    from each voxel under different attention conditions. We then repeated the multivoxel stimulus

584    reconstruction analysis on this simulated data to model population results for the second layer of

585    the model. This approach allowed us to perform virtual experiments to test how changes in the

586    first layer impacted the second layer. That is, we manipulated which vRF parameters changed

587    with attention (first layer) and observed the resulting changes in the population-based stimulus

588    reconstructions (second layer). For example, we could test whether an overall increase in vRF

589    response gain with attention would be necessary or sufficient to reproduce the amplitude

590    increases observed in the empirical stimulus reconstructions reported above. These virtual

591    experiments also allowed us to compare the relative impact of one type of response modulation

592    (e.g. size changes) with other types of response modulations (e.g. position shifts).

593    We first validated our layered model to ensure that it produced results that matched the

594    empirical data. Since our procedure only allowed us to use voxels with reliable vRF fits, we

595    compared the results of the layered model to stimulus reconstructions in the full empirical dataset

596    reported above. **Table 4** reports the RMSE between the fits to the layered IEM reconstructions

597    and the full empirical IEM reconstructions. Our model reproduced the main pattern of results we

598    observed in the previous section. In particular, covert attention led to an increase in the

599    amplitude of reconstructions near the locus of attention, suggesting that our simulated data sets

600    accurately captured the main modulation observed in the real data (**Fig 6b**, gray bars).

601         We then compared two basic manipulations of the layered IEM to see how they

602    contributed to the amplitude increase that we observed in the stimulus reconstructions. When we

603    abolished the position shift between attention conditions in the first layer of the model, we

604    observed a decrease in stimulus reconstruction amplitude (**Fig 6b**). This suggests that spatial

605    position shifts at the single voxel level are necessary for amplitude changes in stimulus

606    reconstructions. When we held vRF sizes constant across attention conditions, there was little

607    change in the amplitude effect in most ROIs, suggesting that size changes in vRFs were not

608    necessary for changes in stimulus reconstruction amplitude.

609         To more formally compare each manipulation of the layered IEM, we compared each

610    model to a baseline in which no vRFs changed with attention (far left in **Fig 6c**). This baseline

611    should have the highest RMSE, and any additional attentional modulations to the underlying

612    vRFs should decrease the error between the simulated data and the empirical data. Conversely, a

613    model with higher RMSE is worse at accounting for the empirical data. In all ROIs, a model that

614    abolished position shifts had a higher RMSE than a model which abolished size shifts (**Fig 6c**,

615    light red and green bars). In fact, just modeling vRF position shifts was sufficient to significantly

616    decrease RMSE in all ROIs except V4. However, this is likely because the layered IEM was a

617    poor model for the attention effects in V4. This is evidenced by the fact that the baseline model

618    did not have the highest RMSE (**Fig 6c**).

619            The overall pattern of results across ROIs is consistent with the interpretation that shifts

620    in the position of vRFs have the largest impact on the population-level representations, while

621    changes in vRF size or gain play smaller roles in changing the fidelity of the population code.

622

623    **DISCUSSION (Max 1500 words; current 1437)**

624    By simultaneously measuring the response properties of both single voxels and populations of

625    voxels within retinotopic areas of visual cortex, we were able to link attentional modulations of

626    spatial encoding properties across scales. Our data provide an initial account of how different

627    types of RF modulations improve the quality of spatial population codes. We first report how

628    different types of vRF modulations depended on the distance between the vRF's preferred

629    position and the static attention target (**Fig 2**). We then found that shifts in the preferred position

630    of vRFs near the attended target increased the spatial discriminability of a population (**Fig 3**), as

631    well as the amplitude of stimulus reconstructions based on populations of vRF responses (**Fig 5**).

632    **Attentional modulations of spatial RFs**

633            While our study is not the first to simultaneously measure a population of RFs tiling a

634    continuous portion of visual space, we provide new data on how vRF responses are modulated

635    around a covertly attended static target (Sprague and Serences, 2013; Klein et al., 2014; Kay et

636    al., 2015; Sheremata and Silver, 2015). Like prior macaque studies, we find that the spatial

637    pattern of attentional modulations is widely variable, but that position shifts in RFs depend

638    heavily on the distance from the attended target (Connor et al., 1996, 1997). We also found that

639    vRF size modulations weakly depended on the RF's spatial relationship to the attended target,

640    even though population-averaged size changes remained at a constant level across voxels (**Fig**

641    **2d-e**). Comparison to the existing literature suggests that patterns of RF size modulations likely

642    depend on the nature of the spatial attention task. In fMRI tasks where subjects attended to the

643    mapping stimulus itself, researchers report average vRF size increases with attention (Sprague

644    and Serences, 2013; Kay et al., 2015; Sheremata and Silver, 2015). Furthermore, the relative size

645    of the attention target and mapping stimulus likely play a key role as well. In macaques, RFs in

646    area MT shrink when measured with a mapping probe smaller than the stimulus, but increase in

647    size when macaques track the mapping probes as they move across the screen (Womelsdorf et

648    al., 2006, 2008; Anton-Erxleben et al., 2009; Niebergall et al., 2011). Taken together, these

649    observations demonstrate that the pattern of response modulations in single cells and in single

650    voxels depends on the spatial relationship between the attended target and the spatial extent of

651    the encoding unit.

652         We note that while the similarity between attentional modulations of single cell RFs and

653    single voxel RFs is compelling, they are not interchangeable. Given that fMRI voxels in

654    retinotopically organized regions of visual cortex sample from a broad array of neurons with

655    roughly the same spatial tuning preferences, a position shift in a vRF could be driven by either a

656    change in the preferred position of single neurons, or by a change in the gain profile across

657    neurons tuned to slightly different locations in the visual field. This principle also holds true of

658    position shifts in single neuron RFs, since cortical neurons typically receive input from those

659    with smaller RFs in earlier visual areas (McAdams and Maunsell, 1999; Baruch and Yeshurun,

660    2014; Dhruv and Carandini, 2014). However, our population-level analyses do demonstrate how

661    the properties of local encoding units–both single-cell and single-voxel RFs–might contribute to

662    population-level representations of space.

663    **Attention boosts the spatial encoding fidelity of a population**

664         We first measured the overall capacity of a population of voxels to make fine spatial

665    discriminations in a region of space. We found that attention increased spatial discriminability

666    near the attended target, relative to the ignored target. We then performed virtual experiments on

667    the underlying vRFs contributing to the population to determine how they affected the spatial

668    discriminability metric. We report that vRF position shifts increased spatial discriminability

669    significantly more than vRF size changes or even gain changes (**Fig 3**). As above, we note that

670    spatially-specific patterns of gain changes in input RFs could produce these position shifts in

671    downstream neural populations. This observation suggests that gain modulations with attention

672    may exert their largest effects on the downstream population, where these patterns of gain

673    changes become apparent shifts in vRF position. Our data are consistent with the interpretation

674    that a neural population only begins to encode the attended area with higher fidelity after input

675    gain changes are transformed into apparent shifts in spatial tuning in the encoding units of that

676    population.

677         Since the spatial discriminability metric is only informative about a local portion of

678    space, we performed a second population analysis to reconstruct an image of the entire visual

679    field on each trial using a multivariate IEM. Attention increased the amplitude of stimulus

680    reconstructions near the attention target, indicating an increase in representational fidelity that

681    accompanies the change in spatial discriminability. In addition, the layered spatial encoding

682    model revealed that shifts in vRF position could sufficiently account for these attentional

683    enhancements in the population-level stimulus reconstructions, but changes in vRF size could

684    not. Altogether, our data demonstrate that spatial tuning shifts in a group of RFs may be the

685    dominant way that single encoding units alter the properties of a population spatial code.

686          Our findings also underscore the fact that changes in the spatial encoding properties of

687     single units do not directly translate into analogous changes in the encoding properties of a

688     population of those same units. This is particularly true when considering the effects of spatial

689     attention on representations of the entire visual scene. Although we found that single units

690     shifted their preferred position towards the attended target, population-level representations did

691     not generally shift with attention. When the population code did shift its encoded position for a

692     given stimulus, we found that it was typically representations of the ignored stimulus that shifted

693     farther from the true stimulus location (**Fig 5**), consistent with a tendency towards more error-

694     prone representations of stimuli far from the relevant portion of the screen. These types of

695     differences further emphasize the need to understand the effects of cognitive state and task

696     demands on population codes for the entire visual scene, rather than focusing solely on single

697     units encoding largely local visual information.

698     **Tuning shifts and labeled lines**

699          Historically, shifts in the tuning of a RF have not been considered one of the main

700     mechanisms by which attention modulates population-level information, although a handful of

701     recent papers suggest that this view is being reconsidered (David et al., 2008; Anton-Erxleben

702     and Carrasco, 2013). This is largely due to 'labeled-line' theories of visual information

703     processing, which posits that a single neuron has a consistent feature label which downstream

704     neurons rely on to perform computations and transmit stable information  (Barlow, 1972;

705     Doetsch, 2000; David et al., 2008). When a spatial RF shifts position as a function of cognitive

706     state (e.g., attention), that single neuron's feature label is no longer consistent. Without an

707     accompanying shift in the downstream neurons receiving the changing feature label, such a

708     change could disrupt the stability of the population code. However, our results suggest that

709    population-level spatial representations remain stable even as the tuning of the underlying vRFs

710    is shifting. In fact, spatial representations are even enhanced as a result of RF shifts.

711         An alternate proposal to a labeled line code relies on the joint information encoded across

712    a population of cells (Erickson, 1982; Doetsch, 2000). This may occur at several scales–for

713    example, V2 could use the pattern of information from V1 inputs to form a visual representation.

714    This idea is more akin to an encoder-decoder model in which the downstream decoder does not

715    need information about the altered representations in each of the encoder units, but instead relies

716    on a population readout rule (Seriès et al., 2009). The population readout rule could incorporate

717    knowledge about the 'labels' of the encoder units, but could perform equally well by relying on

718    relative changes in the pattern across units to resolve the information encoded in the population.

719    This may be a more parsimonious account of the attentional data reported so far. However,

720    further exploration of population readout rules in visual cortex are needed to test this hypothesis.

721

722    **Conclusions**

723         The spatial encoding properties of the visual system can be measured and modeled at

724    many different spatial scales. Here, we report these properties and how they change with

725    attention for single voxels and for a group of voxels in a retinotopic region. Future lines of

726    research into how attention modifies the specific inputs or outputs to a single encoding unit or a

727    population of encoding units may help resolve the question of how shifts in RF labels are

728    generated. Moreover, further investigation into population code readout rules may help

729    adjudicate theories of sensory information processing beyond simple 'labeled line' coding

730    schemes.

731

## REFERENCES

732   **REFERENCES**

733   Amano K, Wandell B a, Dumoulin SO (2009) Visual field maps, population receptive field sizes,

734      and visual field coverage in the human MT+ complex. J Neurophysiol 102:2704–2718.

735   Anton-Erxleben K, Carrasco M (2013) Attentional enhancement of spatial resolution: linking

736      behavioural and neurophysiological evidence. Nat Rev Neurosci 14:188–200.

737   Anton-Erxleben K, Stephan VM, Treue S (2009) Attention reshapes center-surround receptive

738      field structure in macaque cortical area MT. Cereb Cortex 19:2466–2478.

739   Barlow HB (1972) Single units and sensation: A neuron doctrine for perceptual psychology?

740      Perception 1:371–394.

741   Baruch O, Yeshurun Y (2014) Attentional attraction of receptive fields can explain spatial and

742      temporal effects of attention. Vis cogn 22:704–736.

743   Benjamini Y, Yekutieli D (2001) The control of the false discovery rate in multiple testing under

744      dependency. Ann Stat 29:1165–1188.

745   Brouwer GJ, Heeger DJ (2009) Decoding and reconstructing color from responses in human

746      visual cortex. J Neurosci 29:13992–14003.

747   Brouwer GJ, Heeger DJ (2013) Categorical clustering of the neural representation of color. J

748      Neurosci 33:15454–15465.

749   Butts DA, Goldman MS (2006) Tuning Curves, Neuronal Variability, and Sensory Coding. PLoS

750      Biol 4:e92.

751   Cohen MR, Maunsell JHR (2014) Neuronal mechanisms of spatial attention in visual cerebral

752      cortex. In: The Oxford Handbook of Attention, pp 318–345.

753   Connor CE, Gallant JL, Preddie DC, Van Essen DC (1996) Responses in Area V4 Depend on the

754      Spatial Relationship Between Stimulus and Attention. J Neurophysiol 75:1306–1308.

755    Connor CE, Preddie DC, Gallant JL, Van Essen DC (1997) Spatial attention effects in macaque

756        area V4. J Neurosci 17:3201–3214.

757    David S V, Hayden BY, Mazer JA, Gallant JL (2008) Attention to stimulus features shifts

758        spectral tuning of V4 neurons during natural vision. Neuron 59:509–521.

759    Desimone R, Duncan J (1995) Neural mechanisms of selective visual attention. Annu Rev

760        Neurosci 18:193–222.

761    Dhruv NT, Carandini M (2014) Cascaded Effects of Spatial Adaptation in the Early Visual

762        System. Neuron 81:529–535.

763    Doetsch GS (2000) Patterns in the brain. Neuronal population coding in the somatosensory

764        system. Physiol Behav 69:187–201.

765    Dumoulin SO, Wandell BA (2008) Population receptive field estimates in human visual cortex.

766        Neuroimage 39:647–660.

767    Erickson RP (1982) The across-fiber pattern theory: An organizing principle for molar neural

768        function. In: Contributions to sensory physiology, Vol. 6, pp 79–110.

769    Ester EF, Anderson DE, Serences JT, Awh E (2013) A Neural Measure of Precision in Visual

770        Working Memory. J Cogn Neurosci 25:754–761.

771    Ester EF, Sprague TC, Serences JT (2015) Parietal and Frontal Cortex Encode Stimulus-Specific

772        Mnemonic Representations during Visual Working Memory. Neuron 87:893–905.

773    Gattass R, Nascimento-Silva S, Soares JGM, Lima B, Jansen AK, Diogo ACM, Farias MF,

774        Botelho MMEP, Mariani OS, Azzi J, Fiorani M (2005) Cortical visual areas in monkeys:

775        location, topography, connections, columns, plasticity and cortical dynamics. Philos Trans

776        R Soc Lond B Biol Sci 360:709–731.

777    Harvey BM, Dumoulin SO (2011) The relationship between cortical magnification factor and

778    population receptive field size in human visual cortex: constancies in cortical architecture. J

779    Neurosci 31:13604–13612.

780  Huth AG, Nishimoto S, Vu AT, Gallant JL (2012) A continuous semantic space describes the

781    representation of thousands of object and action categories across the human brain. Neuron

782    76:1210–1224.

783  Kay KN, Weiner KS, Grill-Spector K (2015) Attention reduces spatial uncertainty in human

784    ventral temporal cortex. Curr Biol 25:595–600.

785  Klein BP, Harvey BM, Dumoulin SO (2014) Attraction of Position Preference by Spatial

786    Attention throughout Human Visual Cortex. Neuron:1–11.

787  Lee S, Papanikolaou A, Logothetis NK, Smirnakis SM, Keliris G a (2013) A new method for

788    estimating population receptive field topography in visual cortex. Neuroimage 81:144–157.

789  McAdams CJ, Maunsell JHR (1999) Effects of attention on orientation-tuning functions of single

790    neurons in macaque cortical area V4. J Neurosci 19:431–441.

791  Moran J, Desimone R (1985) Selective attention gates visual processing in the extrastriate

792    cortex. Science (80- ) 229:782–784.

793  Navalpakkam V, Itti L (2007) Search Goal Tunes Visual Features Optimally. Neuron 53:605–

794    617.

795  Niebergall R, Khayat PS, Treue S, Martinez-Trujillo JC (2011) Expansion of MT neurons

796    excitatory receptive fields during covert attentive tracking. J Neurosci 31:15499–15510.

797  Pouget A, Dayan P, Zemel RS (2003) Inference and computation with population codes. Annu

798    Rev Neurosci 26:381–410.

799  Regan D, Beverley KI (1985) Postadaptation orientation discrimination. J Opt Soc Am 2:147–

800    155.

801    Roberts M, Delicato LS, Herrero J, Gieselmann M a, Thiele A (2007) Attention alters spatial

802        integration in macaque V1 in an eccentricity-dependent manner. Nat Neurosci 10:1483–

803        1491.

804    Scolari M, Serences JT (2009) Adaptive Allocation of Attentional Gain. J Neurosci 29:11933–

805        11942.

806    Scolari M, Serences JT (2010) Basing perceptual decisions on the most informative sensory

807        neurons. J Neurophysiol 104:2266–2273.

808    Serences JT, Saproo S (2012) Computational advances towards linking BOLD and behavior.

809        Neuropsychologia 50:435–446.

810    Seriès P, Stocker A a, Simoncelli EP (2009) Is the homunculus "aware" of sensory adaptation?

811        Neural Comput 21:3271–3304.

812    Sheremata SL, Silver MA (2015) Hemisphere-Dependent Attentional Modulation of Human

813        Parietal Visual Field Representations. J Neurosci 35:508–517.

814    Sprague TC, Ester EF, Serences JT (2014) Reconstructions of information in visual spatial

815        working memory degrade with memory load. Curr Biol 24:2174–2180.

816    Sprague TC, Ester EF, Serences JT (2016) Restoring Latent Visual Working Memory

817        Representations in Human Cortex. Neuron 91:694–707.

818    Sprague TC, Saproo S, Serences JT (2015) Visual attention mitigates information loss in small-

819        and large-scale neural codes. Trends Cogn Sci:1–12.

820    Sprague TC, Serences JT (2013) Attention modulates spatial priority maps in the human

821        occipital, parietal and frontal cortices. Nat Neurosci 16:1879–1887.

822    Tolias AS, Moore T, Smirnakis SM, Tehovnik EJ, Siapas AG, Schiller PH (2001) Eye

823        movements modulate visual receptive fields of V4 neurons. Neuron 29:757–767.

824    Wandell BA, Winawer J (2015) Computational neuroimaging and population receptive fields.

825        Trends Cogn Sci:1–9.

826    Womelsdorf T, Anton-Erxleben K, Pieper F, Treue S (2006) Dynamic shifts of visual receptive

827        fields in cortical area MT by spatial attention. Nat Neurosci 9:1156–1160.

828    Womelsdorf T, Anton-Erxleben K, Treue S (2008) Receptive field shift and shrinkage in

829        macaque middle temporal area through attentional gain modulation. J Neurosci 28:8934–

830        8944.

831    Zirnsak M, Steinmetz N a, Noudoost B, Xu KZ, Moore T (2014) Visual space is compressed in

832        prefrontal cortex before eye movements. Nature 507:504–507.

833

834

835     **LEGENDS**

836     **Figure 1**. Covert spatial attention task and hypothesized representation changes with shifts of

837     spatial attention. (**a**) Subjects fixated centrally and attended to brief rotations in the pentagon

838     stimulus on the left or right while a flickering checkerboard probe stimulus appeared at one of 51

839     grid locations across the visual field. On control runs, subjects attended to a contrast change at

840     fixation. fMRI data measured during this attention task is used to create visualizable estimates of

841     voxel receptive fields (vRFs) and stimulus reconstructions. (**b**) A receptive field model is fit to

842     the responses of each voxel, and can be described by its x and y position (center), its response

843     baseline, response amplitude, and its size (full-width half maximum). (**c**) Given a population of

844     voxels in a retinotopic region, such as V1, we examine two different measures of spatial

845     information in the population. The first, a spatial discriminability metric, scales with the slope of

846     the tuning curve at a given location in space (**Materials and Methods**). The second relies on a

847     multivariate inverted encoding model (IEM) for space. By reconstructing images of the mapping

848     stimulus on each test trial, we can measure how population-level spatial information changes

849     with attention. We then can model how changes in individual vRFs affect both of these

850     population measures.

851

852     **Figure 2**. Changes in voxel receptive fields (vRFs) across attention conditions. We separately

853     estimated vRFs for every voxel in visual and posterior parietal areas, discarding poorly estimated

854     or noisy voxels (**Table 1**). Unless otherwise specified, figure data is averaged across subjects and

855     error bars show 95% confidence intervals computed with resampling the data distribution. (**a**) An

856     example vRF shows that attending covertly to the left location shifts the center of the receptive

857     field profile to the left, when compared to the neutral attend fixation condition. Voxel is from

858    subject AR in area V3A/B. (**b**) Our vRF estimates reproduced the canonical size-eccentricity

859    relationship (positive slope in all ROIs, p < minimum possible p-value, e.g., 1/1000 iterations)

860    and the increase in slope between visual regions. (**c**) Preferred position changes of V4 vRFs with

861    covert spatial attention. We binned each vRF by its position during the attend fixation condition.

862    The origin of each arrow is the center of each position bin. The end of the arrow shows the

863    average position shift of the vRFs within that position bin during the attend peripheral conditions

864    (left/right are collapsed and shown as attend left). The majority of vRFs shift toward the attended

865    location (blue-green color map vs. red-yellow). (**d**) Mean changes in vRF parameters (attend

866    peripheral target – attend fixation) in each visual area. (**e**) Attentional modulations of each vRF

867    parameter plotted by the vRF's distance from the attention target. We only show areas where

868    these data are significantly described by a polynomial of order n > 0 (**Table 2**).

869

870    **Figure 3**. Spatial discriminability increases with attention and is mediated by position changes in

871    vRFs. Error bars depict bootstrapped 95% CIs. (**a**) We formulated a measurement to describe the

872    ability of a population of voxels to make fine spatial discriminations around the attention target.

873    We used the properties of each voxel's spatial tuning curve to make this measurement

874    (**Materials and Methods**). Spatial discriminability increased when subjects attended the target,

875    compared to when they ignored the target in the opposite hemifield (resampled p < minimum

876    possible p-value (1/1000) for all ROIs for all ROIs). (**b**) The discriminability metric was

877    recomputed for vRFs with a variety of attentional modulations. (*none* = vRF parameters during

878    the neural attend fixation condition; *a* = amplitude; *s* = size; *p* = position). Spatial

879    discriminability increased significantly when we applied position changes measured during the

880    attend L/R task to the vRFs compared to when we applied no parameter changes (solid bar). By

881    contrast, applying size changes either did not change spatial discriminability (V3A/B, V4, IPS0)

882    or even decreased it from the no change condition (V1-V3).

883

884    **Figure 4**. Multivariate inverted encoding model (IEM) used to reconstruct the mapping probe

885    stimuli. (**a**) To train the IEM, we first take the BOLD data from all voxels within a visual region

886    from a subset of training trials. Then, we solve for a set of channel weights using least squares

887    regression. To reconstruct the stimulus, we invert this weight matrix and multiply it with BOLD

888    data from the same voxels during a test trial. This yields a reconstructed channel response

889    profile, which can be transformed into a reconstruction of the mapping stimulus on every trial in

890    each attention condition. Data shown are examples from participant AR for a subset of V1

891    voxels. (**b**) Example stimulus reconstructions for participant AI, V1. These reconstructions were

892    averaged across trials with the same position, yielding 51 reconstructions – one for each unique

893    position in the test dataset. In the left panel, the same averaged position reconstructions are

894    shown for each condition. The amplitude on the left is higher when attending left, and on the

895    right when attending right. (**c**) Average reconstruction sizes and amplitudes for each stimulus

896    position (collapsed across condition; left is attended). The diameter of the circle depicts the

897    average fit FWHM of the reconstructions at that spatial position. Reconstruction amplitude was

898    greater in the attended hemifield compared to the ignored hemifield in areas V3, V4, V3A/B, and

899    IPS0, $p < 0.005$).

900

901    **Figure 5**. Reconstruction parameters as a function of mapping stimulus distance from the

902    covertly attended locations (*s_dist_attn*) and attention hemifield (attended vs. ignored). See

903    **Table 3** for complete list of p-values.

904  **Figure 6**. A layered spatial encoding model reveals how different sets of vRF changes lead to

905  enhancements in multivariate stimulus reconstructions. (**a**) The first layer of the model uses the

906  vRF fits to generate BOLD data from every subject's real trial sequence. Then the BOLD data

907  from all voxels within one ROI is used to train a multivariate spatial encoding model and

908  reconstruct the mapping stimuli. (**b**) Change in reconstruction amplitude in the attended vs. the

909  ignored hemifield. Simulated reconstructions (black) qualitatively reproduce changes in

910  reconstructions using real BOLD data from the same reduced set of voxels (gray bars).

911  Furthermore, while position shifts in vRFs are necessary to observe increases in reconstruction

912  amplitude near the attended location (blue), size changes are not (red). The results for the best

913  model in each visual area is shown in yellow. (**c**) RMSE between each set of IEM fits and the

914  full empirical dataset fits shown in **Fig 5**. The null baseline model (far left) is a layered IEM

915  where the vRF parameters are the same across all attention conditions. We then added vRF

916  attentional modulations for each parameter as shown in the matrix, with the full model on the far

917  right. * indicate an FDR-corrected p-value <.05 for models that differed significantly from the

918  null baseline model. The red bar highlights that a size change model generally performed

919  significantly worse than the null. The green bar highlights that a position change model generally

920  performed significantly better than the null.
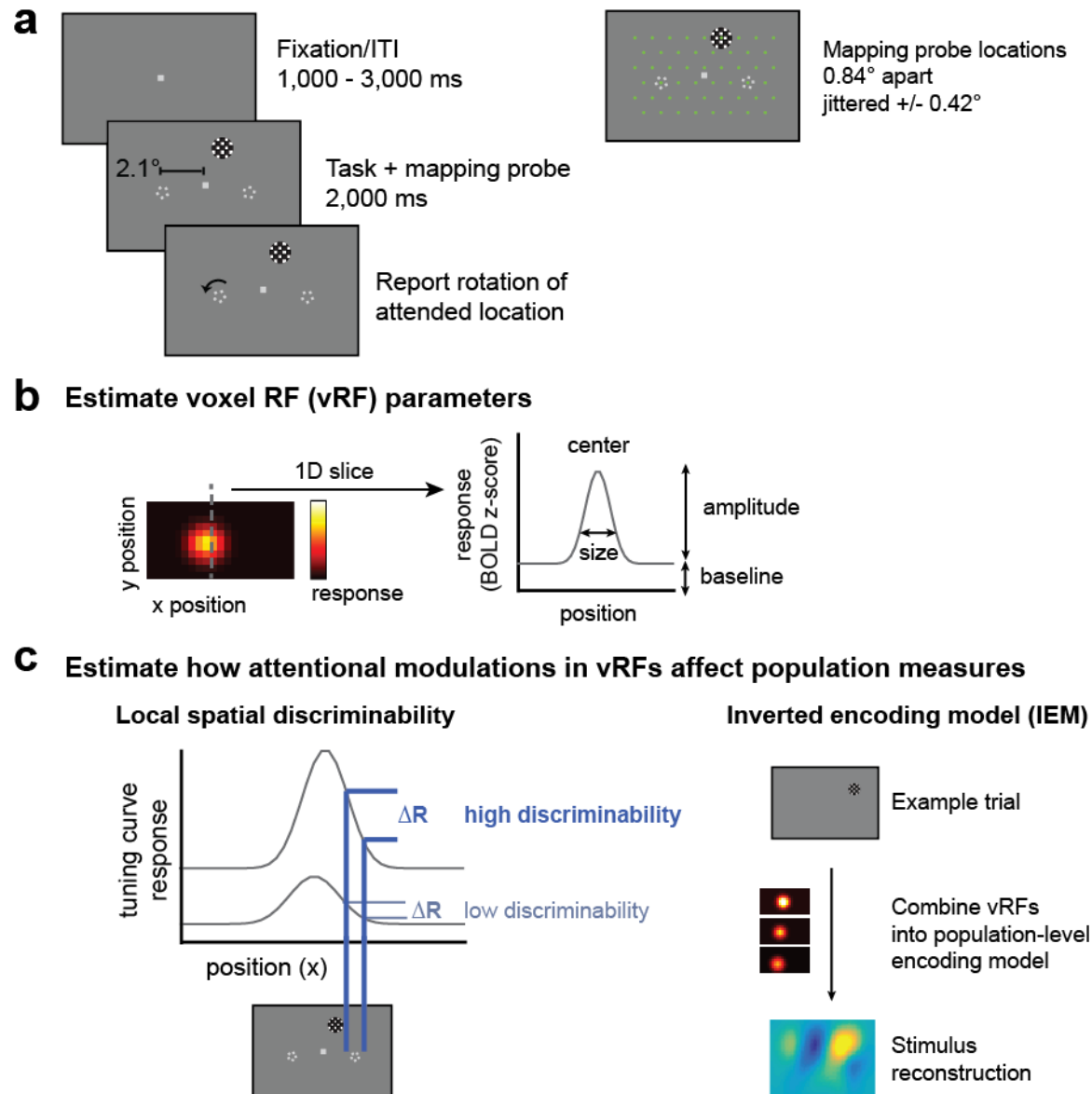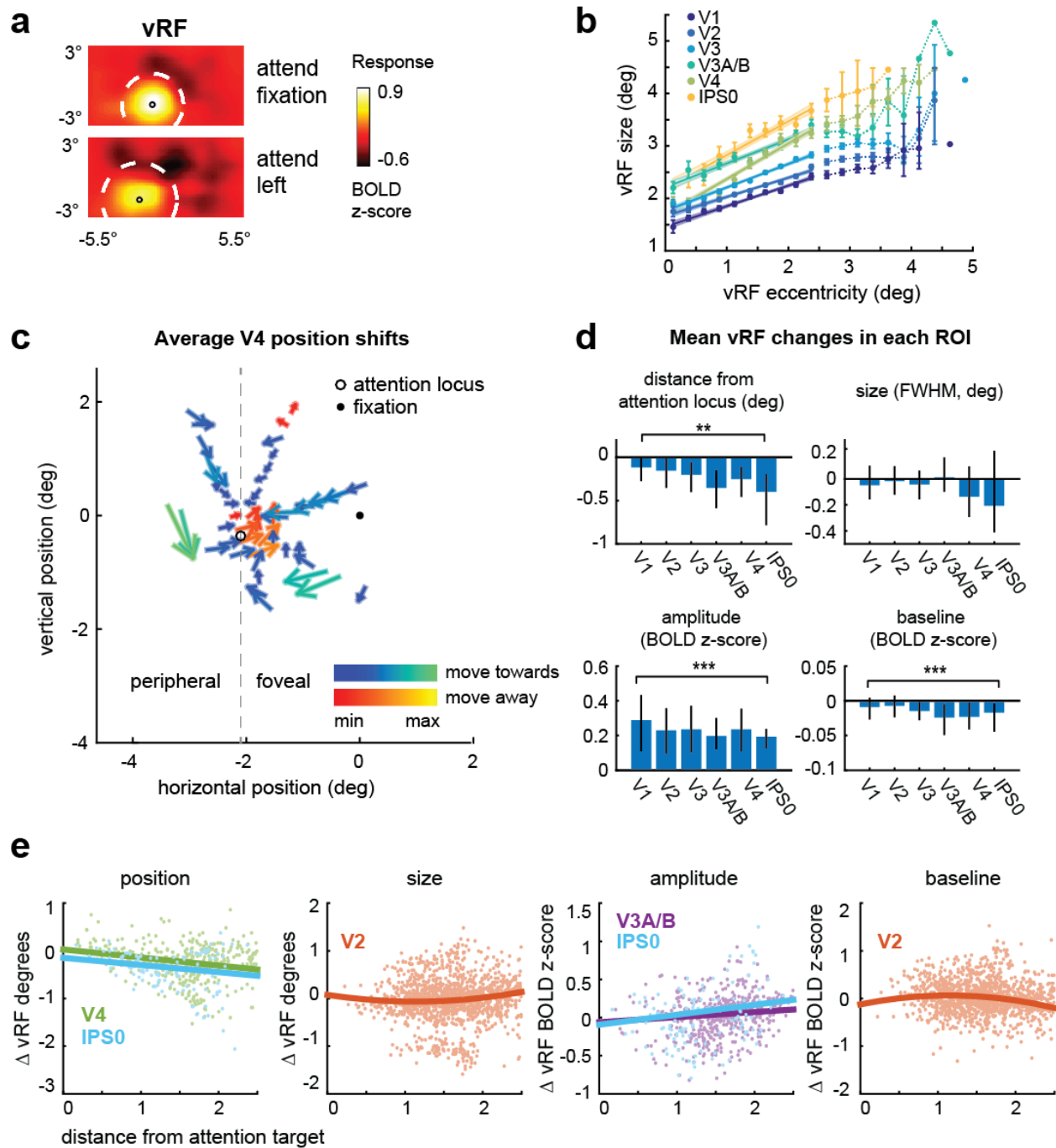
921 **FIGURES**

922 **Figure 1**



923

924

925    **Figure 2**



926    distance from attention target

927

928 **Figure 3**



929

930    **Figure 4**



931

932

933 **Figure 5**

934



935 mapping stimulus distance from attention targets (*s_dist_attn*)
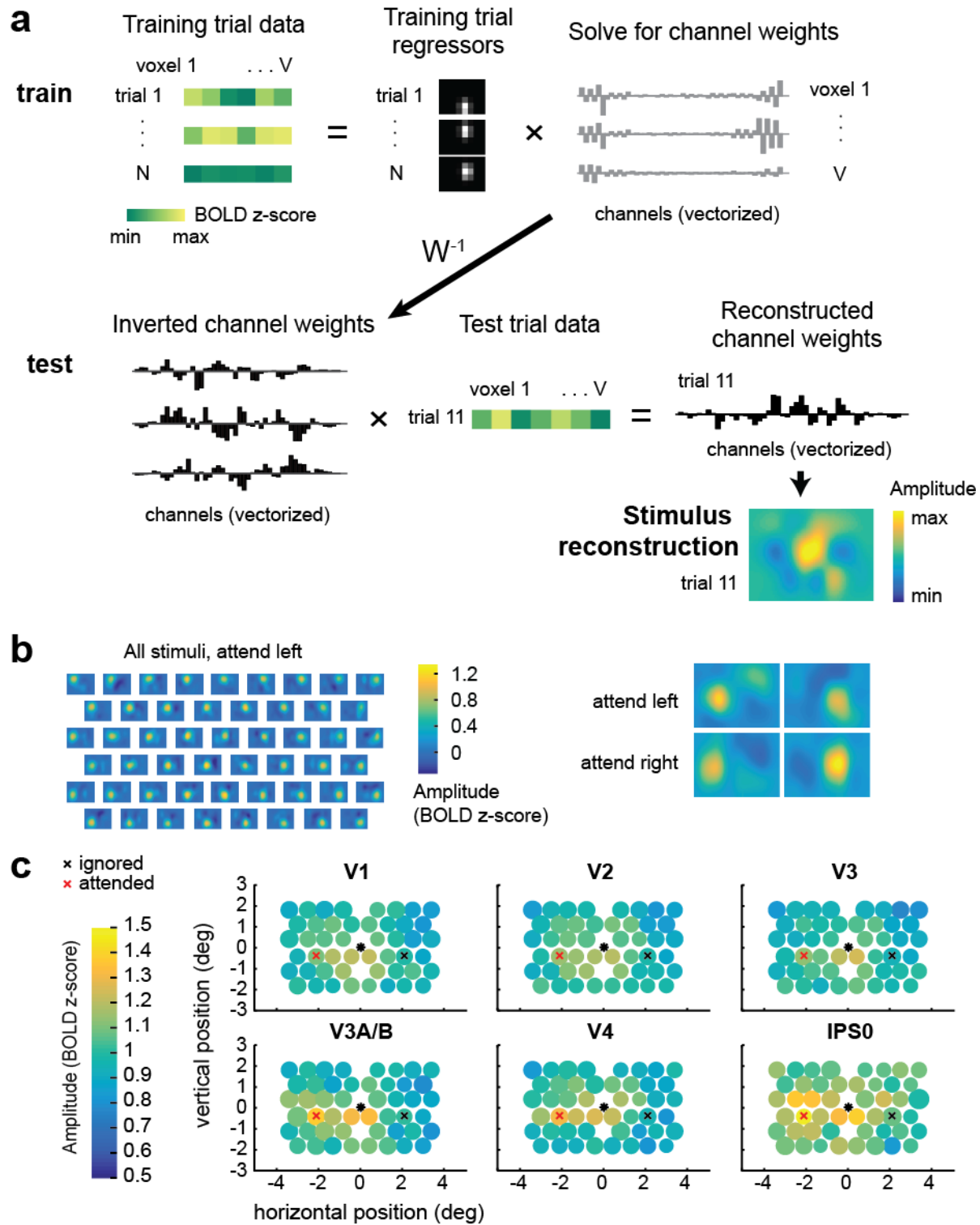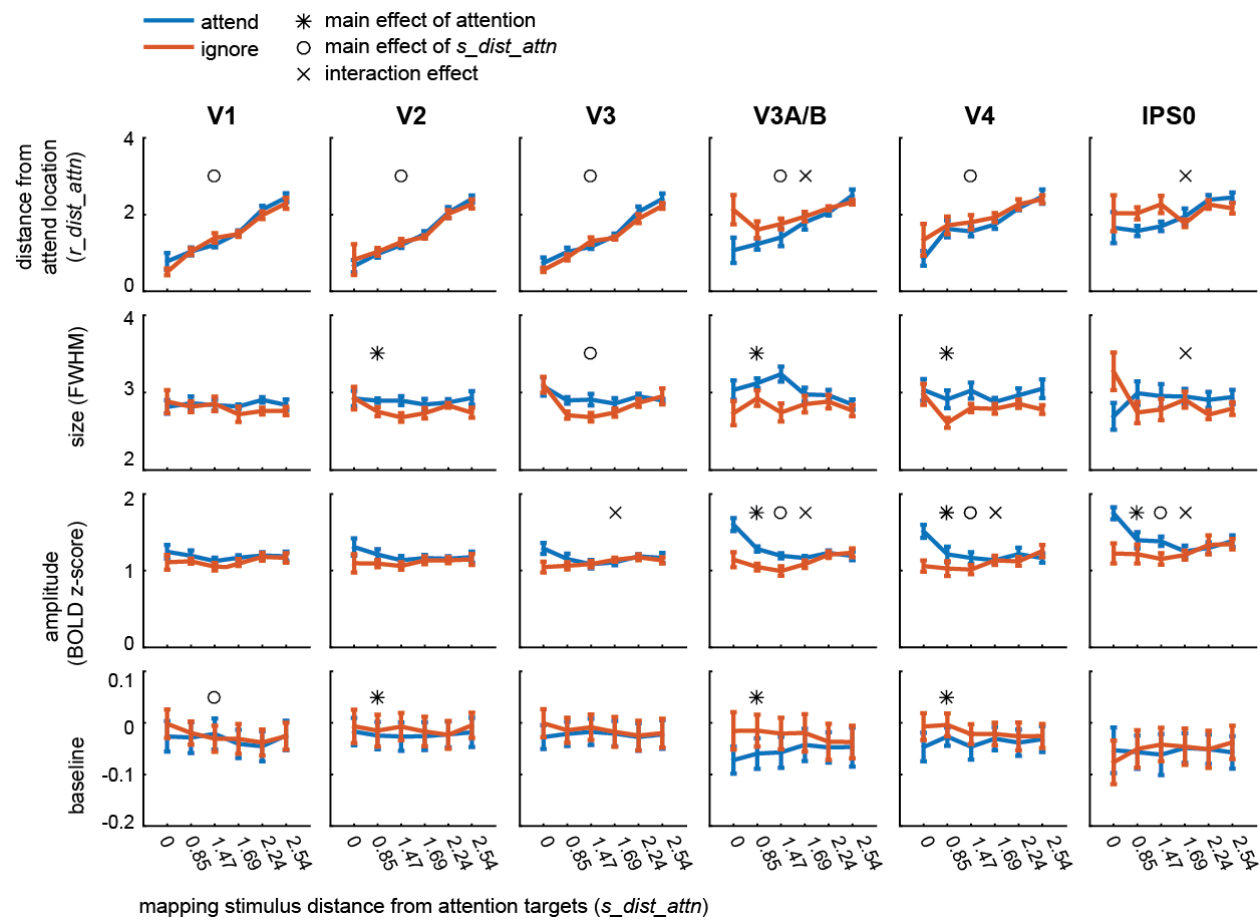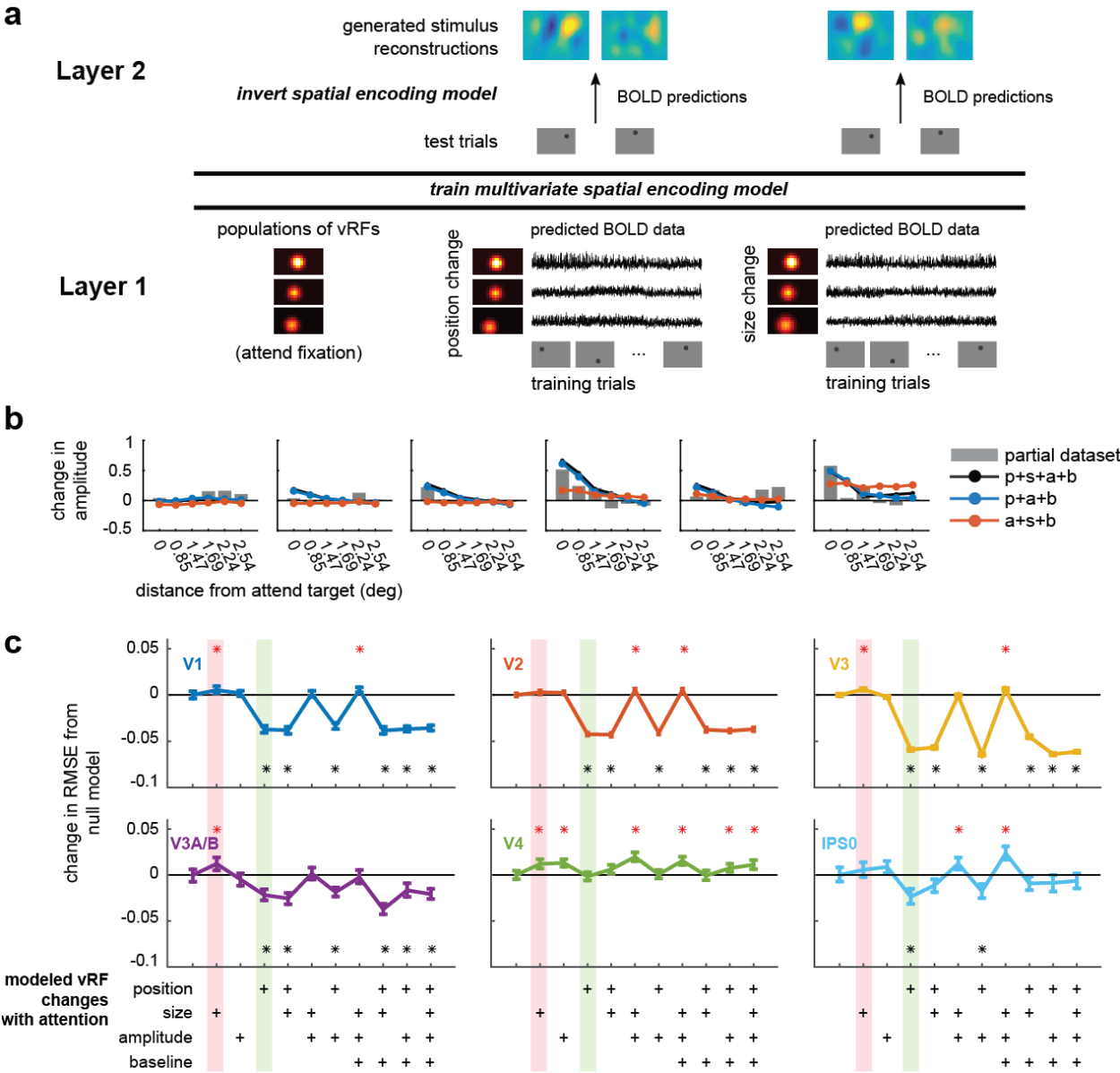
936

937    **Figure 6**



938

939

940    **TABLES**

941    **Table 1**. vRF selection statistics, pooled across participants (N = 7)

| Region of interest | Total number of localized voxels | Number of voxels after GLM thresholding | Number of voxels after regular- izability threshold | Number of voxels after cross- validation threshold | Percent that survive all thresholds | RMSE fit error for surviving voxels |
|---|---|---|---|---|---|---|
| V1 | 3,723 | 3,286 | 2,148 | 969 | 26.03% | 0.1231 |
| V2 | 4,154 | 3,685 | 2,895 | 1,355 | 32.62% | 0.1228 |
| V3 | 3,698 | 3,246 | 2,600 | 1,460 | 39.48% | 0.1191 |
| V3A/B | 1,988 | 1,796 | 1,278 | 446 | 22.43% | 0.1063 |
| V4 | 1,702 | 1,308 | 954 | 349 | 20.51% | 0.1060 |
| IPS0 | 1,430 | 1,250 | 680 | 114 | 7.97% | 0.0985 |
| TOTAL | 16,695 | 14,571 | 10,555 | 4,693 | 28.11% | 0.1126 |

942

943 **Table 2**. Mean coefficients for polynomial fits of how vRF parameter change is modulated by

944 distance from the attended location (*v_dist_attn*)

|  | Position | Size | Amplitude | Baseline |
|---|---|---|---|---|
| V1 | .003 | -.005, .032, -.023 | -.015, -.003 | .039, -.065 |
| V2 | -.022 | *.109*, -.247, .075 | -.031, .076, -.043 | -.*143*, .*328*, -.123 |
| V3 | **-.103** | .007, -.015, -.025 | .011, -.037 | .128, -.221 |
| V3A/B | **-.332** | -.*012*, -.017, .063 | *.067*, -.061 | -.122 |
| V4 | **-.168**, .028 | -.056, .138, -.043 | .002 | -.039, .105, -.181 |
| IPS0 | -.*151*, **-.146** | -.080, .291, -.192 | **.131, -.094** | -.*269* |

945

946 [a] **bold** numbers indicate that the p-value passed FDR-correction (q = 0.05) across ROIs and

947 coefficients within each parameter; *italicized* numbers are $p < .05$, uncorrected. Number of

948 reported coefficients in the table correspond to the polynomial order which was yielded the most

949 parsimonious fit to the data (e.g., 1 coefficient for n = 0, 2 coefficients for n = 1, etc.).

950

951

952    **Table 3**. 2-way ANOVA results for reconstruction parameter changes (*s_dist_attn* x attention

953    hemifield).

| | V1 | V2 | V3 | V3A/B | V4 | IPS0 |
|---|---|---|---|---|---|---|
| Omnibus test | | | | | | |
|     Position | **<.001** | **<.001** | **<.001** | **<.001** | **<.001** | **.007** |
|     Size | .917 | .097 | **.001** | **.001** | **.017** | **.016** |
|     Amplitude | .207 | .220 | **.003** | **<.001** | **<.001** | **<.001** |
|     Baseline | .024 | .257 | .485 | **.002** | **.004** | .925 |
| Main effect of distance | | | | | | |
|     Position | **<.001** | **<.001** | **<.001** | **<.001** | **<.001** | .084 |
|     Size | | | **.001** | .233 | .169 | .679 |
|     Amplitude | | | .269 | **<.001** | **.008** | **.007** |
|     Baseline | | | | .864 | .336 | |
| Main effect of attention | | | | | | |
|     Position | .573 | .920 | .399 | .022 | .189 | .235 |
|     Size | | | .163 | **.001** | **.005** | .509 |
|     Amplitude | | | .047 | **.005** | **.002** | **.028** |
|     Baseline | | | | **<.001** | **.001** | |
| Interaction of distance & attention | | | | | | |
|     Position | .188 | .892 | .354 | **.001** | .679 | **.004** |
|     Size | | | .157 | .099 | .582 | **.005** |
|     Amplitude | | | **.003** | **<.001** | **<.001** | **.002** |
|     Baseline | | | | .210 | .202 | |

954

955    [a] **bold** numbers indicate that the p-value passed FDR-correction (q = .05, corrected across ROIs

956    and comparisons within each parameter).

957 **Table 4**. RMSE (and 95% CIs) between reconstructions from the reduced dataset (only using voxels with RFs) or from different

958 versions of the layered IEM using the same voxels.

| | Real data | p/s/a/b | p/a/b | p/s/b | s/a/b | p/a | s/a | p/s | p | a | s | none |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| V1 | 0.773 [0.640, 0.913] | 0.184 [0.181, 0.187] | 0.183 [0.180, 0.186] | 0.181 [0.177, 0.185] | 0.225 [0.222, 0.228] | 0.186 [0.183, 0.189] | 0.221 [0.218, 0.224] | 0.181 [0.177, 0.185] | 0.182 [0.178, 0.186] | 0.221 [0.218, 0.224] | 0.225 [0.221, 0.229] | 0.219 [0.216, 0.223] |
| V2 | 0.394 [0.326, 0.461] | 0.185 [0.183, 0.187] | 0.184 [0.182, 0.185] | 0.185 [0.182, 0.187] | 0.227 [0.226, 0.229] | 0.181 [0.180, 0.183] | 0.228 [0.226, 0.229] | 0.179 [0.177, 0.182] | 0.180 [0.178, 0.182] | 0.225 [0.223, 0.226] | 0.225 [0.223, 0.227] | 0.222 [0.221, 0.224] |
| V3 | 0.368 [0.280, 0.480] | 0.180 [0.179, 0.181] | 0.177 [0.176, 0.179] | 0.196 [0.194, 0.198] | 0.247 [0.246, 0.249] | 0.177 [0.176, 0.178] | 0.241 [0.239, 0.242] | 0.184 [0.183, 0.186] | 0.182 [0.181, 0.184] | 0.239 [0.238, 0.241] | 0.247 [0.246, 0.249] | 0.241 [0.240, 0.243] |
| V3A/B | 0.745 [0.522, 0.978] | 0.336 [0.331, 0.342] | 0.340 [0.333, 0.348] | 0.320 [0.314, 0.326] | 0.354 [0.348, 0.362] | 0.338 [0.334, 0.343] | 0.359 [0.352, 0.365] | 0.331 [0.325, 0.337] | 0.335 [0.329, 0.342] | 0.352 [0.345, 0.359] | 0.369 [0.362, 0.376] | 0.357 [0.350, 0.363] |
| V4 | 0.729 [0.617, 0.843] | 0.393 [0.389, 0.398] | 0.390 [0.386, 0.394] | 0.381 [0.376, 0.387] | 0.398 [0.393, 0.402] | 0.383 [0.378, 0.388] | 0.402 [0.397, 0.407] | 0.388 [0.384, 0.393] | 0.381 [0.376, 0.385] | 0.395 [0.391, 0.399] | 0.394 [0.389, 0.399] | 0.382 [0.378, 0.387] |
| IPS0 | 0.879 [0.764, 1.006] | 0.491 [0.483, 0.499] | 0.489 [0.480, 0.498] | 0.489 [0.481, 0.496] | 0.521 [0.514, 0.529] | 0.480 [0.473, 0.488] | 0.510 [0.504, 0.517] | 0.486 [0.479, 0.493] | 0.474 [0.466, 0.483] | 0.506 [0.500, 0.513] | 0.503 [0.495, 0.512] | 0.498 [0.491, 0.506] |

959

960 [a] To generate CIs, the resampling of the real data is performed at the level of the fits to the reconstructions, whereas resampling

961 layered IEM RMSEs is described in **Materials and Methods**