

1 **Bottom-up and top-down computations in word- and face-selective cortex**

2

3 Kendrick N. Kay <sup>\*1</sup> & Jason D. Yeatman <sup>\*2</sup>

4

5 <sup>1</sup>Center for Magnetic Resonance Research, Department of Radiology, University of Minnesota, Twin  
6 Cities, Minneapolis, MN, 55455, USA

7 <sup>2</sup>Institute for Learning & Brain Sciences and Department of Speech & Hearing Sciences, University of  
8 Washington, Seattle, WA, 98195, USA

9 \*Corresponding authors ([kay@umn.edu](mailto:kay@umn.edu), [jyeatman@uw.edu](mailto:jyeatman@uw.edu))

10

11 **ABSTRACT**

12

13 The ability to read a page of text or recognize a person's face depends on category-selective visual regions  
14 in ventral temporal cortex (VTC). To understand how these regions mediate word and face recognition, it  
15 is necessary to characterize how stimuli are represented and how this representation is used in the  
16 execution of a cognitive task. Here, we show that the response of a category-selective region in VTC can  
17 be computed as the degree to which the low-level properties of the stimulus match a category template.  
18 Moreover, we show that during execution of a task, the bottom-up representation is scaled by the  
19 intraparietal sulcus (IPS), and that the level of IPS engagement reflects the cognitive demands of the task.  
20 These results provide an account of neural processing in VTC in the form of a model that addresses both  
21 bottom-up and top-down effects and quantitatively predicts VTC responses.

22

23 **INTRODUCTION**

24

25 How does visual cortex work? One approach to answering this question consists in building functional  
26 models that characterize the computations that are implemented by neurons and their circuitry (Hubel and  
27 Wiesel, 1963; Heeger et al., 1996). This approach has been fruitful for the front end of the visual system,  
28 where relatively simple image computations have been shown to characterize the spiking activity of  
29 neurons in the retina, thalamus, and V1 (Carandini et al., 2005; Wu et al., 2006). Based on this pioneering  
30 work in electrophysiology, researchers have extended the modeling approach to characterize responses in  
31 human visual cortex, as measured by functional magnetic resonance imaging (fMRI) (Wandell, 1999;  
32 Dumoulin and Wandell, 2008; Kay et al., 2008).

33

34 Models of early visual processing have been able to offer accurate explanations of low-level perceptual  
35 functions such as contrast detection (Ress et al., 2000; Ress and Heeger, 2003) and orientation  
36 discrimination (Bejjanki et al., 2011). However, these models are insufficient to explain high-level  
37 perceptual functions such as the ability to read a page of text or recognize a face. These abilities are  
38 believed to depend on category-selective regions in ventral temporal cortex (VTC), but the computations  
39 that give rise to category-selective responses are poorly understood.

40  
41 The goal of the present study is to develop a model that predicts fMRI responses in high-level visual  
42 cortex of human observers while they perform different cognitive tasks on a wide range of images. We  
43 seek a model that is fully computable—that is, a model that can operate on any arbitrary visual image and  
44 quantitatively predict BOLD responses and behavior (Kay et al., 2008; Huth et al., 2012; Khaligh-Razavi  
45 and Kriegeskorte, 2014; Yamins et al., 2014). Achieving this goal requires four innovations: First, we  
46 need to develop a forward model that characterizes the relationship between visual inputs and the BOLD  
47 response in word- and face-selective cortex. Second, we need to dissociate bottom-up stimulus-driven  
48 effects from modulation by top-down cognitive processes and characterize how these processes alter the  
49 stimulus representation. Third, we need to localize the source of the top-down effects and integrate  
50 bottom-up and top-down computations into a single consolidated model. Finally, the neural computations  
51 should be linked to the measured behavior of the visual observer. In this study, we make progress on these  
52 four innovations and develop a model that characterizes bottom-up and top-down computations in word-  
53 and face-selective cortex.

## 54 55 **RESULTS**

### 56 57 **VTC responses depend on both stimulus properties and cognitive task**

58  
59 Ventral temporal cortex (VTC) is divided into a mosaic of high-level visual regions that respond  
60 selectively to specific image categories, and are believed to play an essential role in object perception  
61 (Kanwisher, 2010; Dehaene and Cohen, 2011). We focus on two specific VTC regions, the visual word  
62 form area (VWFA), which selectively responds to words (Cohen et al., 2000; 2002; Wandell et al., 2012),  
63 and the fusiform face area (FFA), which selectively responds to faces (Kanwisher et al., 1997; Grill-  
64 Spector and Weiner, 2014).

65  
66 We measured blood oxygenation level dependent (BOLD) responses to a set of carefully controlled  
67 images while manipulating the cognitive task that the subjects performed on the stimuli. The first task

68 was designed to minimize the influence of cognitive processes on sensory processing of the stimulus.  
69 Subjects performed a demanding perceptual task on a small dot ( $0.12^\circ \times 0.12^\circ$ ) presented at fixation. In  
70 this *fixation task*, the presented stimuli are irrelevant to the subject, and we interpret evoked activity as  
71 reflecting primarily the intrinsic, bottom-up response from VTC. We acknowledge that the fixation task  
72 may not perfectly isolate bottom-up responses. For example, high-contrast stimuli may automatically  
73 attract attention. Moreover, there are other potential interpretations of the fixation task: for example,  
74 allocating attention to the small fixation dot might engage active suppression of responses to the  
75 presented stimuli.

76  
77 To a first approximation, much of the variance in the bottom-up fixation responses from VWFA and FFA  
78 is explained by the category of stimulus (**Figure 1d, red lines**). However, we find that responses are not  
79 invariant to low-level properties of the stimulus: both image contrast and phase coherence modulate  
80 response amplitudes. For example, the response to a word in VWFA is 2.4 times stronger when the word  
81 is presented at 100% contrast as compared to 3% contrast. These bottom-up effects (see also Rainer et al.,  
82 2001; Avidan et al., 2002; Yue et al., 2011; Nasr et al., 2014) may be somewhat surprising given that  
83 theories of word recognition generally posit that the VWFA response is invariant to low-level features  
84 (Dehaene and Cohen, 2007; 2011; Price and Devlin, 2011). In fact, it is currently debated whether the  
85 VWFA should be considered a visual area or a “meta modal” language region (Reich et al., 2011; Striem-  
86 Amit et al., 2012). Our measurements indicate that when top-down signals are minimized, word- and  
87 face-selective cortex is sensitive to low-level image properties, and that an accurate model of the  
88 computations performed by these regions must consider not only the stimulus category but also low-level  
89 features of the stimulus.

90

91 \*\*\* Insert Figure 1 here \*\*\*

92

93 We also measured VTC responses while subjects performed a *categorization task*, in which the subject  
94 reports the perceived category of the stimulus, and a *one-back task*, in which the subject detects  
95 consecutive repetitions of stimulus frames. Despite presentation of identical stimuli across the three tasks,  
96 there are substantial changes in evoked VTC responses. Responses are larger for the categorization  
97 (**Figure 1d, green lines**) and one-back tasks (**Figure 1d, blue lines**) compared to the fixation task  
98 (**Figure 1d, red lines**), and we interpret these response increases as reflecting top-down modulation. In  
99 some cases, the top-down modulation is even larger than the modulation achieved by manipulation of the  
100 stimulus. For example, the VWFA response to 3%-contrast words during the one-back task exceeds the  
101 response to 100%-contrast words. Note that the task effects cannot be explained simply by differences in

102 spatial attention: the one-back task produces substantially larger responses than the categorization task  
103 despite the fact that both tasks require the locus of spatial attention to be on the stimulus. Task effects in  
104 lower-level areas exist but are smaller in size (**Figure 1–figure supplement 1**).

105

106 A potential explanation of the top-down modulation is differences in task difficulty (Ress et al., 2000).  
107 For example, it is presumably more difficult to perceive low-contrast stimuli than high-contrast stimuli,  
108 and this may explain why there is large response enhancement for low- but not high-contrast stimuli (see  
109 VWFA contrast-response function for word stimuli in **Figure 1d**). Later in this paper, we provide a  
110 computational mechanism that could underlie the psychological concept of task difficulty.

111

112 In summary, our measurements indicate that VTC responses cannot be interpreted without specifying the  
113 cognitive state of the observer. A complete model of the computations performed by VWFA and FFA  
114 must consider the cognitive task in addition to stimulus properties.

115

### 116 **Model of bottom-up computations in VTC**

117

118 Before addressing the influence of top-down factors, we first develop a model of bottom-up responses in  
119 VWFA and FFA. Although the field has long understood that stimulus category is a good predictor of  
120 evoked responses (Kanwisher et al., 1997; Kriegeskorte et al., 2008; Grill-Spector and Weiner, 2014), we  
121 do not yet have a computational explanation of this phenomenon. In other words, although we are able to  
122 use our own visual systems to assign a label such as “word” or “face” to describe the data, we have not  
123 yet identified the operations that enable our visual systems to derive these labels in the first place. An  
124 additional limitation of our conceptual understanding is that it fails to account for the sensitivity of  
125 VWFA and FFA to low-level image properties. We therefore ask: Is it possible to develop a quantitative  
126 characterization of the bottom-up computations that can reproduce observed stimulus selectivity in human  
127 VTC?

128

129 Extending an existing computational model of fMRI responses in the visual system (Kay et al., 2008;  
130 2013b; 2013c), we conceive of a model involving two stages of image computations (**Figure 2a**). The  
131 first stage consists of a set of local oriented filters, akin to what has been used to model physiological  
132 responses in V1 (Jones and Palmer, 1987; Carandini et al., 2005). The second stage consists of a  
133 normalized dot product applied to the outputs of the first stage. This dot product computes how well a  
134 given stimulus matches a category template (e.g., a word template for VWFA, a face template for FFA).  
135 We construct category templates directly from the stimulus set used in the experiment; in a later section

136 we explore how well this approach generalizes. The present model, termed the *Template model*, is almost  
137 certainly an oversimplification of the complex nonlinear processing performed in VTC. Nevertheless, the  
138 model is theoretically motivated, consistent with hierarchical theories of visual processing (Fukushima,  
139 1980; Heeger et al., 1996; Serre et al., 2007; DiCarlo et al., 2012; Rolls, 2012), and provides a useful  
140 starting point for characterizing the computations that underlie word- and face-selectivity. Moreover,  
141 unlike recently popular deep neural network models that also involve hierarchical processing (Khaligh-  
142 Razavi and Kriegeskorte, 2014; Yamins et al., 2014; Güçlü and van Gerven, 2015), the model we propose  
143 is parsimonious with only three free parameters, and is therefore straightforward to fit and interpret (see  
144 **Figure 2a** and Methods).

145  
146 \*\*\* Insert Figure 2 here \*\*\*

147  
148 Applying the Template model to responses measured during the fixation task, we find that the Template  
149 model accurately predicts a large amount of variance in the responses of VWFA and FFA (**Figure 2b**).  
150 The model outperforms a phenomenological model, termed the *Category model*, that posits that perceived  
151 stimulus category is sufficient to predict the response of category-selective regions. Notably, the  
152 Template model is able to predict the response to non-preferred stimulus categories in each ROI. This  
153 suggests that responses to non-preferred stimuli are meaningful and the result of a well-defined  
154 computation performed by the visual system (Haxby et al., 2001). The model also outperforms simplified  
155 versions of the Template model that include only one of the two processing stages, as well as versions of  
156 the Template model in which the category template lacks tuning (non-selective template), is equally  
157 weighted between words and faces (mixed template), or is constructed randomly (random template)  
158 (**Figure 2c**).

159  
160 The experiment we have conducted explores a limited range of stimuli. To further assess how well the  
161 Template model generalizes, we collected an additional dataset that includes 92 images taken from a  
162 previous study of object representation (Kriegeskorte et al., 2008). In its original instantiation (**Figure**  
163 **2a**), the Template model uses category templates tailored to the stimuli in the main experiment, and we  
164 find that this instantiation of the Template model does not generalize well to the larger stimulus set  
165 (**Figure 2–figure supplement 1b**). However, by implementing a simple model extension in which we use  
166 a data-driven approach to estimate category templates, we find that the Template model achieves a  
167 reasonable level of accuracy on the new stimulus set (**Figure 2–figure supplement 1d**). This finding  
168 validates the basic architecture of the Template model, demonstrates how the Template model might be  
169 extended to account for increasingly large ranges of measurements, and provides a promising method to

170 model response properties of other high-level visual regions not investigated here (e.g. place- and limb-  
171 selective cortex).

172

173 The Template model advances us towards a computational understanding of VTC by demonstrating that  
174 BOLD responses in VTC can be predicted based on a template-matching operation on incoming visual  
175 inputs filtered by early visual cortex. The present results indicate that although high-level representations  
176 are not identical to low-level properties, they are built from, and fundamentally tied to, low-level  
177 properties through a series of linear and nonlinear operations. This conclusion is consistent with classic  
178 hierarchical theories of visual cortex (Fukushima, 1980; Heeger et al., 1996; Serre et al., 2007; DiCarlo et  
179 al., 2012; Rolls, 2012) and recent evidence that visual features may explain semantic representations  
180 found in high-level visual cortex (Jozwik et al., 2016). Our model can be viewed as a potential  
181 mechanism for how semantic tuning properties emerge in visual cortex (Huth et al., 2012). Our results  
182 indicate when studying high-level sensory representations in the brain, a precise characterization of the  
183 stimulus still matters.

184

### 185 **Top-down modulation acts as a stimulus-specific scaling**

186

187 While the Template model explains bottom-up responses of VWFA and FFA as indexed by the fixation  
188 task, it does not explain why responses are higher in these areas during the categorization and one-back  
189 tasks (see **Figure 1d**). This is simply because the stimuli are identical across the three tasks and the  
190 response of the Template model, like that of many computational models of visual processing, is solely a  
191 function of the stimulus. Before we can design a model to capture the top-down effects, we must first  
192 understand exactly how top-down signals shape the VTC response.

193

194 By visualizing VTC responses as points in a multi-dimensional neural space with VWFA, FFA, and hV4  
195 BOLD response amplitudes as the axes, we see that responses to words and faces lie on specific  
196 manifolds, appearing as “arms” that emanate from the origin (**Figure 3**). Importantly, we observe that the  
197 categorization and one-back tasks act as a scaling mechanism on the representation observed during the  
198 fixation task. The scaling mechanism moves the representation of each stimulus along the arms and away  
199 from the origin. Moreover, the amount of scaling is not constant across stimuli but is stimulus-specific,  
200 and this is most evident when considering the lowest contrast stimuli (**Figure 3, black dots**).

201

202

\*\*\* Insert Figure 3 here \*\*\*

203

204 The visualization also shows that substantial responses to non-preferred categories are present in each  
205 ROI (e.g., faces in VWFA, words in FFA) and that these responses are scaled during the stimulus-directed  
206 tasks. Thus, not only is information regarding non-preferred categories present in each ROI, but this  
207 information is actively modulated when subjects perform a perceptual task on those categories. These  
208 observations support the view that the brain uses a distributed strategy for perceptual processing and that  
209 category-selective regions are components of a more general network of regions that coordinate to extract  
210 visual information (Haxby et al., 2001; Cox and Savoy, 2003). An alternative scheme, more in line with a  
211 modular view of perceptual processing (Kanwisher and Wojciulik, 2000; Baldauf and Desimone, 2014),  
212 is area-specific enhancement, in which the representation of a stimulus is enhanced only in the region that  
213 is selective for that stimulus (e.g., enhancement of words only in VWFA, enhancement of faces only in  
214 FFA). This scheme is not supported by our measurements (**Figures 3b and 3c**; formal model evaluation  
215 is performed in a later section). Rather, response scaling occurs even for non-preferred stimulus  
216 categories, and the amount of scaling varies as a function of stimulus properties such as image contrast.

217  
218 A simple interpretation of the scaling effects is that they serve to increase signal-to-noise ratio in visually  
219 evoked responses in VTC (Brouwer and Heeger, 2013). For example, assuming that one use of the  
220 stimulus representation in VTC is to discriminate whether the presented stimulus is a word or face (or,  
221 more generally, identify the category of the stimulus (DiCarlo et al., 2012)), the scaling induced by the  
222 stimulus-directed tasks serves to increase the distance of neural responses from a linear decision boundary  
223 that separates words and faces (**Figure 3d**). Interestingly, the categorization and one-back tasks appear to  
224 act via the same scaling mechanism. The stronger scaling observed for the one-back task might be a  
225 consequence of increased amplitude or duration of neural activity. These results suggest that, at least for  
226 the perceptual tasks sampled here and the spatial scale of neural activity measured in this study, top-down  
227 cognitive processes do not impart additional tuning or selectivity but serve to amplify the selectivity that  
228 is already computed by visual cortex.

### 229 230 **IPS is the source of top-down modulation to VTC**

231  
232 To design a plausible model that can predict top-down effects, we next turn to identifying the neural  
233 circuitry that generates task modulations in VTC. There are two candidate mechanisms. The first is that  
234 sensitivity to task is locally generated from the neuronal architecture of VTC itself. We explore an  
235 alternative hypothesis whereby top-down modulation is induced by input from another brain region that is  
236 sensitive to task demands. To identify this region, we perform a connectivity analysis in which we first  
237 subtract the bottom-up signal in VTC, as given by responses measured during the fixation task, from

238 responses measured during the categorization and one-back tasks. We then correlate these residuals,  
239 which isolate the top-down signal, against the responses of every cortical location.

240

241 Applying this connectivity analysis to our data, we find that responses in the intraparietal sulcus (IPS)  
242 predict the top-down enhancement of VTC responses (**Figure 4b**) better than responses in any other  
243 region of cortex. As a control, if we omit the subtraction step and simply correlate raw VTC responses  
244 with the responses of different cortical locations, we find that the correlation is instead strongest with a  
245 range of areas spanning occipital cortex (**Figure 4a**). This indicates that the VTC response is a mixture of  
246 bottom-up and top-down effects and that the top-down influence from the IPS becomes clear only when  
247 bottom-up effects are removed. Comparing our results to a publicly available atlas (Wang et al., 2014),  
248 we estimate that the source of top-down modulation is localized to the IPS-0 and IPS-1 subdivisions of  
249 the IPS (see also **Figure 4—figure supplement 1**).

250

251 \*\*\* Insert Figure 4 here \*\*\*

252

253 Previous research has identified IPS as playing a key role in controlling spatial attention (Saalmann et al.,  
254 2007; Lauritzen et al., 2009). Our results extend these findings by showing that, despite the fact that  
255 spatial attention is always directed towards the foveal stimulus during the categorization and one-back  
256 tasks, the amount of modulation from the IPS is flexible and varies depending on properties of the  
257 stimulus and demands of the task. For example, during the categorization task, the observed enhancement  
258 for low-contrast stimuli is much larger than that for high-contrast stimuli. This mechanism could explain  
259 the finding that difficult tasks enhance visual responses (Ress et al., 2000).

260

261 The direct influence of IPS on neural responses in VTC is consistent with anatomical measurements  
262 demonstrating the existence of a large white-matter pathway connecting dorsal and ventral visual cortex,  
263 called the vertical occipital fasciculus (VOF) (Yeatman et al., 2013; 2014; Takemura et al., 2016). Using  
264 diffusion-weighted MRI and tractography (data acquired in 8 of 9 subjects), we show that the VOF  
265 specifically connects the VWFA and FFA with the functionally identified peak region in the IPS (**Figure**  
266 **4c**). The VWFA falls within the ventral terminations of the VOF for seven subjects and, for the eighth,  
267 the VWFA is 2.7 mm anterior to the VOF, well within the margin of error for tractography (Jeurissen et  
268 al., 2011). The FFA falls within the ventral terminations of the VOF for all eight subjects. These results  
269 provide an elegant example of how anatomy subserves function, and sets the stage for a circuit-level  
270 computational model that, guided by anatomical constraints, characterizes the computations that emerge  
271 from interactions between multiple brain regions.

272

273 **Model of top-down computations in VTC**

274

275 The previous two sections provide critical insights that set the stage for building a quantitative model that  
276 predicts top-down effects in VTC. Building upon the observation that top-down modulation acts as a  
277 scaling mechanism on responses in VWFA and FFA (see **Figure 3**) and the observation that top-down  
278 effects are correlated with the IPS signal (see **Figure 4**), we propose that the magnitude of the IPS  
279 response to a stimulus indicates the amount of top-down scaling that is applied to bottom-up sensory  
280 responses in VTC (**Figure 5a**). We implement this model, termed the *IPS-scaling model*, using response  
281 magnitudes extracted from a broad anatomical mask of the IPS. This strategy helps avoid the overfitting  
282 that might ensue from a more specific voxel-selection procedure tailored to the fine-scale and potentially  
283 idiosyncratic pattern of results from the connectivity analysis. For example, if we were to select the single  
284 cortical location in the IPS that best correlates with the top-down modulation of VTC, this would make  
285 voxel selection a critical part of the model and render the modeling analysis circular (Kriegeskorte et al.,  
286 2009). Nevertheless, the selection procedure is not completely independent, so the modeling results  
287 should not be viewed as providing independent evidence for the involvement of the IPS.

288

289

\*\*\* Insert Figure 5 here \*\*\*

290

291 We find that the IPS-scaling model accurately characterizes the observed data (**Figure 5b**). For example,  
292 notice that the FFA response to faces increases gradually for each contrast increment during the fixation  
293 task (relatively unsaturated contrast-response function, red arrow). When subjects perform the one-back  
294 task, we observe a U-shaped contrast-response function in IPS (green arrow); multiplication of the two  
295 functions predicts a contrast-response function that is highly saturated and accurately matches the  
296 observed contrast-response function in FFA during the one-back task (blue arrow).

297

298 Importantly, the IPS-scaling model uses a single set of scale and offset parameters on the IPS response  
299 and accurately predicts scaling of VTC responses across the categorization and one-back tasks (**Figure**  
300 **5b, top plot**). This finding suggests that the scaling of VTC by IPS is a general mechanism supporting  
301 perception and is independent of the specific cognitive task performed by the observer. Furthermore, the  
302 scale and offset parameters that are estimated from the data show that when IPS exhibits close to zero  
303 evoked activity (e.g. FACE at 100%-contrast; see **Figure 1–figure supplement 1**), the corresponding  
304 scaling factor is close to one. This has a sensible interpretation: when IPS is inactive, we observe only the  
305 bottom-up response in VTC and no top-down modulation.

306

307 We assessed the cross-validation performance of the IPS-scaling model in comparison to several  
308 alternative models of top-down modulation (including those schematized earlier in **Figure 3b**). In line  
309 with earlier observations (**Figures 3b and 3c**), we find that a model positing enhancement for only the  
310 preferred stimulus category of each area (Area-specific enhancement model) does not optimally describe  
311 the data. We find that a phenomenological scaling model (Scaling model (task-specific)) outperforms a  
312 phenomenological additive model (Additive model (task-specific)), confirming earlier observations that  
313 the top-down modulation is a scaling effect (**Figure 5c**). This conclusion is further supported by the  
314 higher performance observed when the IPS interacts with VTC multiplicatively (IPS-scaling model)  
315 compared to when it interacts additively (IPS-additive model). Finally, we find that the performance of  
316 the IPS-scaling model degrades if the IPS input into the model is shuffled across conditions (IPS-scaling  
317 model (shuffle, shuffle within task)), confirming that top-down modulation from the IPS is dependent on  
318 the stimulus and task.

319

320 Is the IPS the only region that induces top-down modulation of VTC? Inspection of the connectivity  
321 results (see **Figure 4b**) reveals that the top-down residuals in VTC are correlated, to a lesser extent, with  
322 responses of other regions. These weaker correlations might be incidental, or might capture other  
323 important signals. Given that the IPS-scaling model accounts for nearly all of the variance induced by  
324 top-down modulation of VTC (see **Figures 5b and 5c**), we suggest that it is sufficient to consider only  
325 the IPS for the current set of measurements. However, future measurements that employ new stimulus  
326 manipulations and other cognitive tasks may reveal the role of a more extensive brain network. The IPS-  
327 scaling model can be extended to account for new measurements by systematically parameterizing the  
328 connectivity with additional brain regions. For example, some models of reading posit that language-  
329 related regions can directly influence the VWFA (Twomey et al., 2011), suggesting that to account for  
330 measurements made during a more naturalistic reading task, it may be necessary to include Broca's area  
331 in the model.

332

### 333 **Model of perceptual decision-making in IPS**

334

335 Although informative, the finding that IPS provides top-down stimulus-specific scaling of VTC is an  
336 incomplete explanation, as the burden of explaining the top-down effects is simply shifted to the IPS. We  
337 are thus left wondering: is it possible to explain the response profile of the IPS? In particular, can we  
338 explain why the IPS is more active for certain stimuli compared to others? Answering these questions will  
339 provide a critical link between cognitive state and IPS activity.

340

341 Inspired by previous research on perceptual decision-making (Shadlen and Newsome, 2001; Heekeren et  
342 al., 2004; Gold and Shadlen, 2007; Kayser et al., 2010a), we implement a *Drift diffusion model* that  
343 attempts to account for IPS responses measured during the categorization task (**Figure 6a**). The model  
344 uses VTC responses during the fixation task as a measure of sensory evidence, and posits that the IPS  
345 accumulates evidence from VTC over time and exhibits an activity level that is monotonically related to  
346 accumulation time. For example, when VTC responses are small, as is the case for low-contrast stimuli,  
347 sensory evidence for stimulus category is weak, leading to long accumulation times (indexed by  
348 measurements of reaction time during the experiment), and large IPS responses.

349

350

\*\*\* Insert Figure 6 here \*\*\*

351

352 Our implementation of the Drift diffusion model involves two steps. First, we use VTC responses during  
353 the fixation task (reflecting sensory evidence) to predict reaction times measured in the categorization  
354 task. The quality of the predictions is quite high (**Figure 6b, left**). Second, we apply a simple monotonic  
355 function to the reaction times measured during the categorization task to predict the level of response in  
356 the IPS (see Methods). The rationale is that neural activity in IPS is expected to be sustained over the  
357 duration of the decision-making process (Shadlen and Newsome, 2001), and so the total amount of neural  
358 activity integrated over time should be larger for longer decisions. Assuming that the BOLD signal  
359 reflects convolution of a sluggish hemodynamic response function and fine-scale neural activity  
360 dynamics, small differences in the duration of neural activity (e.g. between 0–2 s) are expected to  
361 manifest in differences in BOLD amplitudes (Kayser et al., 2010a) and only minimally in the shapes of  
362 BOLD timecourses. The cross-validated predictions of our proposed model explain substantial variance in  
363 IPS (**Figure 6b, right**).

364

365 It is possible to offer a psychological explanation of IPS activity as reflecting task difficulty—for  
366 example, we can posit that IPS activity is enhanced for low-contrast stimuli because the observer works  
367 harder to perceive these stimuli. The value of the model we have proposed is that it provides a  
368 quantitative and formal explanation of the computations that underlie ‘difficulty’. According to the  
369 model, categorization of low-contrast stimuli is difficult because the IPS computations required to  
370 perform the task involve longer accumulation time, and this is reflected in the fact that IPS response  
371 magnitudes increase monotonically with reaction time. Thus, our model performs several critical  
372 functions: it relates the cognitive task performed by the subject to IPS activity, proposes a computational  
373 explanation of task difficulty, and posits that top-down modulation of VTC by IPS is a direct consequence

374 of fulfilling task demands. We have substantiated this hypothesis for the categorization task and suggest  
375 that this will serve as a foundation for modeling more complex cognitive tasks such as one-back.

376

## 377 **DISCUSSION**

378

379 In summary, we have measured and modeled how bottom-up and top-down factors shape responses in  
380 word- and face-selective cortex. A template operation on low-level visual properties generates a bottom-  
381 up stimulus representation, while top-down modulation from the IPS scales this representation in service  
382 of the behavioral goals of the observer. We develop a computational approach that posits explicit models  
383 of the information processing performed by a network of interacting sensory and cognitive regions of the  
384 brain and validate this model on experimental data. We make publicly available data and open-source  
385 software code implementing the model at <http://cvnlab.net/vtcipsmodel/>.

386

387 The model we propose is valuable because it integrates and explains a range of different stimulus and task  
388 manipulations that affect responses in VWFA, FFA, and IPS. Response properties in these regions can  
389 now be interpreted using a series of simple, well-defined computations that can be applied to arbitrary  
390 images. However, it is also important to recognize the limitations of the model. First, we have tested the  
391 model on only a limited range of stimuli and cognitive tasks. Second, the accuracy with which the model  
392 accounts for the data is reasonable but by no means perfect. For instance, the Template model does not  
393 capture the step-like response profile of VTC as phase coherence is varied (see **Figure 2b**), and the  
394 accuracy with which the model accounts for a wide range of stimuli that includes faces, animals, and  
395 objects, is moderate at best (see **Figure 2–figure supplement 1**). Third, we have thus far characterized  
396 VTC and IPS responses at only a coarse spatiotemporal scale (i.e., BOLD responses averaged over  
397 specific regions-of-interest). Given these limitations, the present work constitutes a first step towards the  
398 goal of developing a comprehensive computational model of human high-level visual cortex. We have  
399 provided data and code so that other researchers can build on our approach, for example, by testing the  
400 generalizability of the model to other stimuli and tasks, extending and improving the model, and  
401 comparing the model against alternative models.

402

403 The fact that cognitive factors substantially affect stimulus representation in visual cortex highlights the  
404 importance of tightly controlling and manipulating cognitive state when investigating stimulus selectivity.  
405 In the present measurements, the most striking example comes from stimulus contrast. When subjects  
406 perform the fixation task, the contrast-response function (CRF) in VWFA is monotonically increasing,  
407 whereas during the one-back task, the CRF flips in sign and is monotonically decreasing (see **Figure 1d**).

408 This effect (similar to what is reported in Murray and He, 2006) is puzzling if we interpret the CRFs as  
409 indicating sensitivity to the stimulus contrast, but is sensible if we interpret the CRFs as instead reflecting  
410 the interaction of stimulus properties and cognitive processes. The influence of cognition on visual  
411 responses forces us to reconsider studies that report unexpected tuning properties in VTC and IPS, such as  
412 tuning to linguistic properties of text in VWFA (Vinckier et al., 2007) and object selectivity in parietal  
413 cortex (Sereno and Maunsell, 1998). In experiments that do not tightly control the cognitive processes  
414 executed by the observer, it is impossible to distinguish sensory effects from cognitive effects. Our  
415 quantitative model of VTC-IPS interactions provides a principled baseline on which to re-interpret past  
416 findings, design follow-up experiments, and guide data analysis.

417  
418 There is a large body of literature on characterizing and modeling the effect of spatial attention on  
419 contrast-response functions (Luck et al., 1997; Boynton, 2009; Reynolds and Heeger, 2009; Itthipuripat et  
420 al., 2014). Although several models have been proposed, none of these straightforwardly account for the  
421 present set of measurements. The *response-gain* model (McAdams and Maunsell, 1999) posits that  
422 attention causes a multiplicative scaling of contrast-response functions. Our observations are consistent  
423 with the general notion of response scaling (see **Figure 3**), but importantly, we find that the amount of  
424 scaling differs for different stimuli. Whereas the response-gain model implies that scaling is constant and  
425 therefore contrast-response functions should grow steeper during the stimulus-directed tasks, we find the  
426 opposite (see **Figure 1d**). The *contrast-gain* model (Reynolds et al., 2000) posits that attention causes a  
427 leftward shift of contrast-response functions (as if contrast were increased). This model does not account  
428 for our measurements, since the stimulus-directed tasks can generate responses to low-contrast stimuli  
429 that are larger than responses to high-contrast stimuli (see **Figure 1d**). Finally, the *additive-shift* (Buracas  
430 and Boynton, 2007) model posits that attention causes an additive increment to contrast-response  
431 functions; we find that our observations are better explained by a scaling, not additive, mechanism (see  
432 **Figures 3 and 5c**). Thus, the effects we report are novel, and previous models of attention cannot explain  
433 these effects. Furthermore, by investigating responses to a wide range of stimuli (including manipulations  
434 of not only contrast but also phase coherence and stimulus category) and by characterizing the source of  
435 attentional signals, our work develops a more comprehensive picture of information processing in the  
436 visual system.

437  
438 There are a number of research questions that remain unresolved. First, our connectivity analysis and  
439 modeling of VTC-IPS interactions is based on correlation of BOLD responses and does not provide  
440 information regarding the directionality, or timing, of neural interactions. In other words, our correlational  
441 results do not, in and of themselves, prove that the IPS causes VTC modulation; rather, we are imposing

442 an interpretation of the results in the context of a computational model. We note that our interpretation is  
443 in line with previous work on perceptual decision-making showing top-down influence (Granger  
444 causality) of IPS on visual cortex (Kayser et al., 2010b). Our working hypothesis is that sensory  
445 information arrives at VTC (as indexed by fixation responses), these signals are routed to IPS for  
446 evidence accumulation, and then feedback from the IPS modulates the VTC response (as indexed by  
447 categorization and one-back responses). Temporally resolved measurements of neural activity (e.g. EEG,  
448 MEG, ECoG) will be necessary to test this hypothesis. Second, the scaling of BOLD response amplitudes  
449 by IPS is consistent with at least two potential mechanisms at neural level: the IPS may be inducing a  
450 scaling on neural activity in VTC or, alternatively, a sustainment of neural activity in VTC. Some support  
451 for the latter comes from a recent study demonstrating that ECoG responses in FFA exhibit sustained  
452 activity that is linked to long reaction times in a face gender discrimination task (Ghuman et al., 2014).  
453 Finally, IPS is part of larger brain networks involved in attention (Corbetta and Shulman, 2002) and  
454 decision-making (Gold and Shadlen, 2007), and identifying the computational roles of other regions in  
455 these networks is necessary for a comprehensive understanding of the neural mechanisms of perception.

456

#### 457 **ACKNOWLEDGMENTS**

458

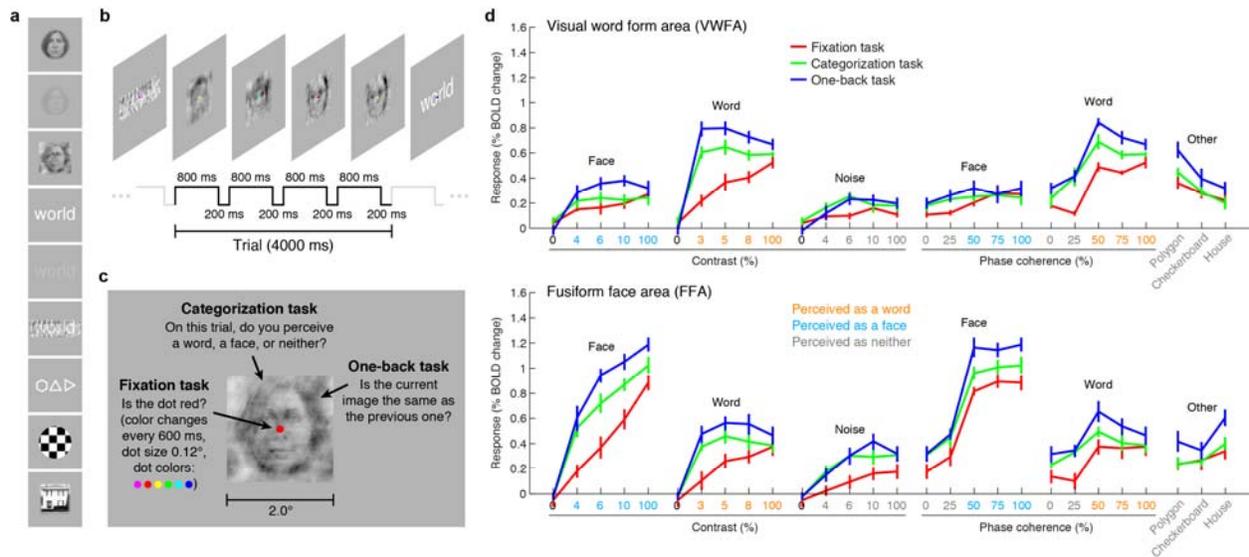
459 We thank K. Grill-Spector for providing the face and house stimuli used in the main experiment, R. Kiani  
460 and N. Kriegeskorte for providing the object stimuli used in the retinotopic mapping experiment, A. Vu  
461 and E. Yacoub for collecting pilot data, C. Gratton, M. Harms, and L. Ramsey for scanning assistance,  
462 and K. Weiner for assistance with ROI definition. We also thank P. Elder, C. Gratton, S. Petersen, A.  
463 Rokem, A. Vogel, and J. Winawer for helpful discussions. This work was supported by the McDonnell  
464 Center for Systems Neuroscience and Arts & Sciences at Washington University (K.N.K.) and NSF Grant  
465 BCS-1551330 (J.D.Y.). Computations were performed using the facilities of the Washington University  
466 Center for High Performance Computing, which were partially provided through grant NCRR  
467 1S10RR022984-01A1.

468

#### 469 **COMPETING INTERESTS**

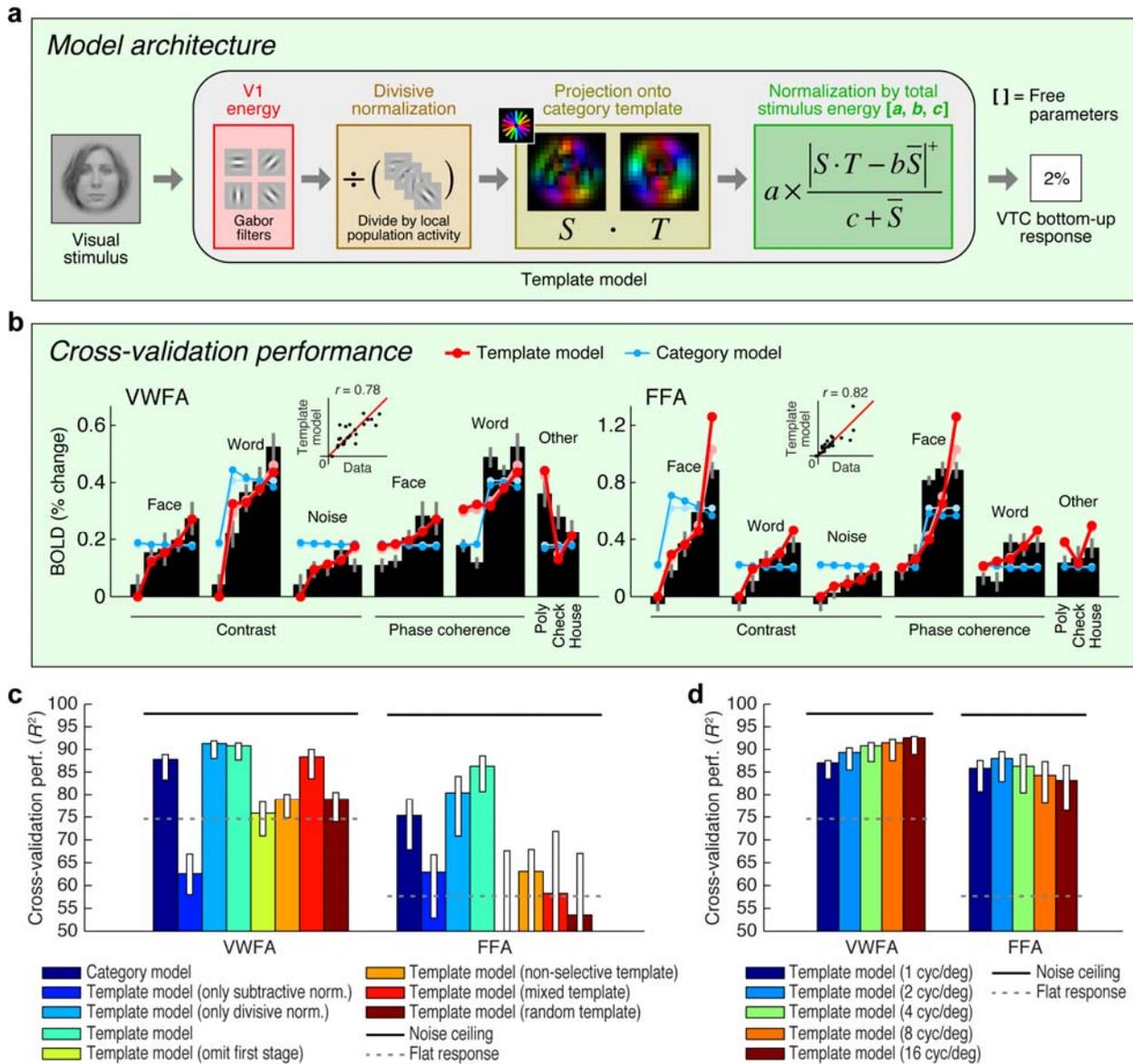
470

471 The authors declare no competing interests.



472  
 473  
 474  
 475  
 476  
 477  
 478  
 479  
 480  
 481

**Figure 1. VTC responses depend on both stimulus properties and cognitive task.** (a) *Stimuli*. Stimuli included faces, words, and noise patterns presented at different contrasts and phase-coherence levels, as well as full-contrast polygons, checkerboards, and houses. (b) *Trial design*. Each trial consisted of four images drawn from the same stimulus type. (c) *Tasks*. On a given trial, subjects performed one of three tasks. (d) *Evoked responses in VWFA (top) and FFA (bottom) for different stimuli and tasks*. Color of x-axis label indicates the perceived stimulus category as reported by the subjects. Error bars indicate bootstrapped 68% CIs.



483

484

485 **Figure 2. Model of bottom-up computations in VTC.** (a) *Model architecture*. The predicted response of

486 the Template model is given by a series of image computations (see Methods). (b) *Cross-validation*

487 *performance*. Black bars indicate bottom-up stimulus-driven responses measured during the fixation task,

488 dark lines and dark dots indicate model predictions (leave-one-stimulus-out cross-validation), and light

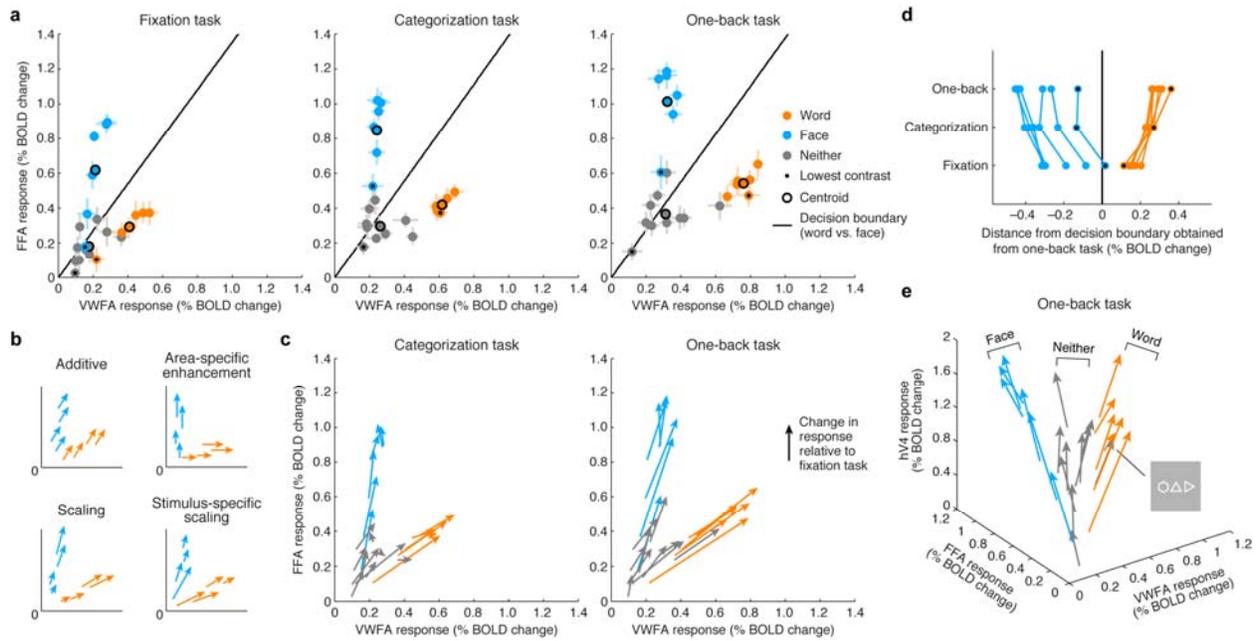
489 lines and light dots indicate model fits (no cross-validation). Scatter plots in the inset compare model

490 predictions against the data. The Template model is compared to the Category model which simply

491 predicts a fixed response level for stimuli from the preferred stimulus category and a different response

492 level for all other stimuli (the slight decrease in response as a function of contrast is a result of the cross-

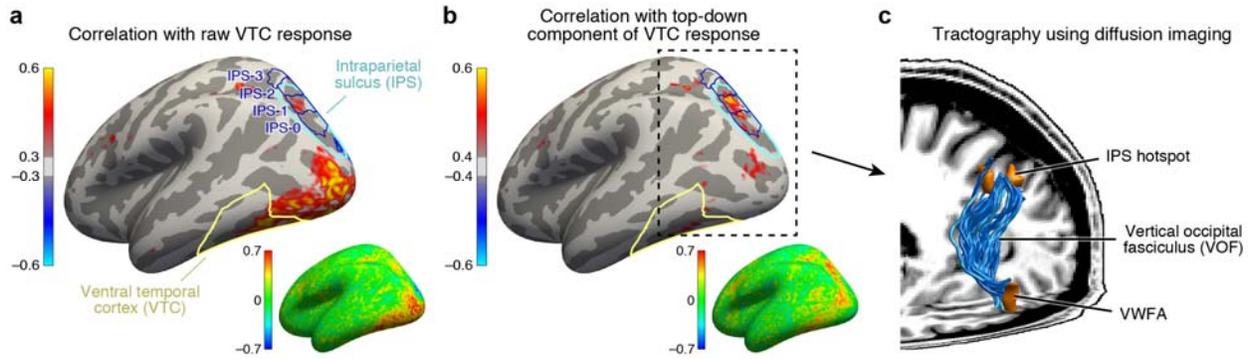
493 validation process). (c) *Comparison of performance against control models*. Bars indicate leave-one-  
494 stimulus-out cross-validation performance. Error bars indicate 68% CIs, obtained by bootstrapping  
495 (resampling subjects with replacement). Solid horizontal lines indicate the noise ceiling, i.e., the  
496 maximum possible performance given measurement variability in the data. Dotted horizontal lines  
497 indicate the cross-validation performance of a model that predicts the same response level for each data  
498 point (this corresponds to  $R^2 = 0$  in the conventional definition of  $R^2$  where variance is computed relative  
499 to the mean). The performance of the Template model degrades if the second stage of nonlinearities is  
500 omitted (Template model (only subtractive normalization)) or if the first stage of the model involving V1-  
501 like filtering is omitted (Template model (omit first stage)). The plot also shows that the precise  
502 configuration of the template is important for achieving high model performance (Template model (non-  
503 selective, mixed, random templates)). (d) *Performance as a function of spatial frequency tuning*. Here we  
504 manipulate the spatial frequency tuning of the filters in the Template model (while fixing spatial  
505 frequency bandwidth at 1 octave). The Template model uses a single set of filters at a spatial frequency  
506 tuning of 4 cycles/degree.  
507



508

509

510 **Figure 3. Top-down stimulus-specific scaling of VTC representation.** (a) Responses plotted in multi-  
 511 dimensional neural space. Each dot indicates ROI (VWFA, FFA) responses to a stimulus. In each plot,  
 512 the black line indicates a linear decision boundary separating words and faces (nearest-centroid classifier,  
 513 angular distance). (b) Schematics of potential top-down mechanisms (these models are formally evaluated  
 514 in **Figure 5c**; see Methods section ‘IPS-scaling model’ for details). (c) Categorization and one-back tasks  
 515 produce stimulus-specific scaling. Arrows indicate the change in representation compared to the fixation  
 516 task. (d) Scaling improves readout. Each data point indicates the signed Euclidean distance between the  
 517 word-face decision boundary (as determined from the one-back task) and the neural response to a single  
 518 stimulus. Lines join data points that correspond to the same stimulus. The scaling observed during the  
 519 categorization and one-back tasks moves responses away from the decision boundary, thereby improving  
 520 signal-to-noise ratio. (e) Separation of other stimulus categories. Including hV4 as a third dimension  
 521 reveals that stimuli categorized as neither words nor faces manifest as a third “arm” that emanates from  
 522 the origin. Although not reported to be a word by the subjects, the polygon stimulus behaves similarly to  
 523 word stimuli.



524

525

526 **Figure 4. IPS is the source of top-down modulation to VTC.** (a) *Correlation with raw VTC response.*

527 This map depicts the correlation between the VTC response observed during the categorization and one-

528 back tasks with the response at each cortical location (inset shows an unsmoothed and unthresholded

529 map). Positive correlations are broadly distributed across occipital cortex. Results are shown for subjects

530 with whole-brain coverage ( $n = 3$ ); results for other subjects with partial-brain coverage ( $n = 6$ ) are shown

531 in **Figure 4–figure supplement 1**. (b) *Correlation with top-down component of VTC response.* After

532 removing bottom-up responses (fixation task), the correlation is spatially localized to a hotspot in IPS-0/1.

533 (c) *Tractography using diffusion MRI.* We find that the vertical occipital fasciculus (Yeatman et al., 2014)

534 connects VWFA and FFA to the IPS hotspot in each subject for which diffusion data were collected ( $n =$

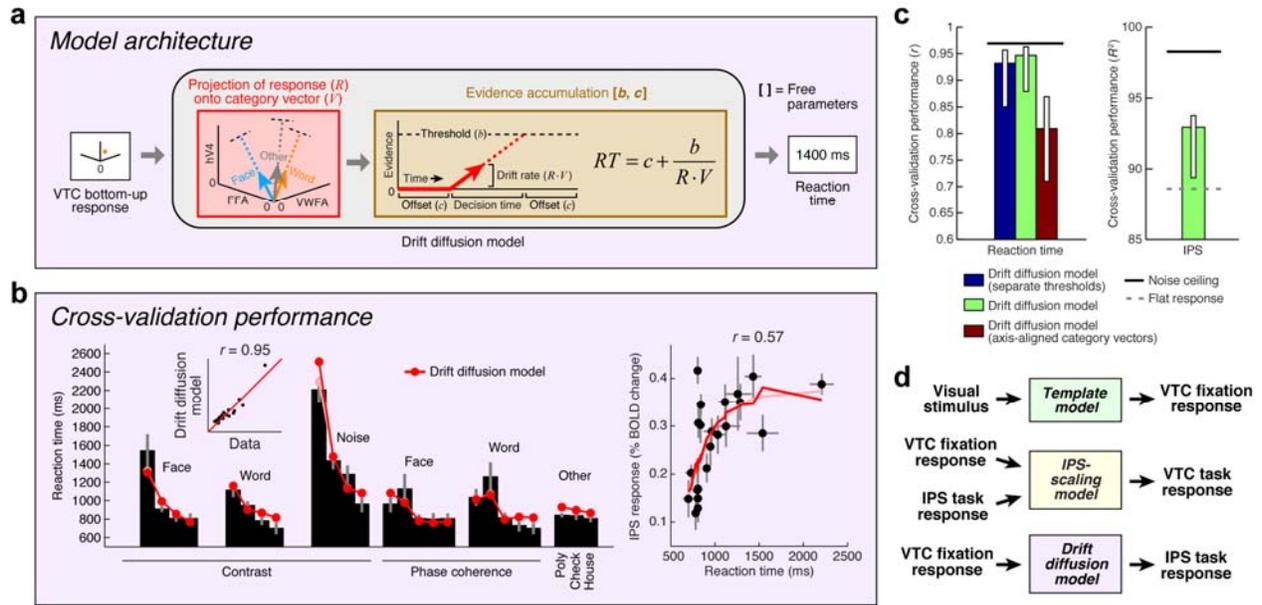
535 8) (rendering shows a representative subject).

536

537

538





553  
 554  
 555  
 556  
 557  
 558  
 559  
 560  
 561  
 562  
 563  
 564  
 565  
 566  
 567  
 568  
 569  
 570  
 571

**Figure 6. Model of perceptual decision-making in IPS.** (a) *Model architecture.* We implement a model that links the stimulus representation in VTC to a decision-making process occurring in IPS. The model first uses the bottom-up VTC response as a measure of sensory evidence and predicts reaction times in the categorization task. The model then predicts the IPS response as a monotonically increasing function of reaction time. Note that this model does not involve stochasticity in the evidence-accumulation process, and is therefore a simplified version of the classic drift diffusion model (Ratcliff, 1978). (b) *Cross-validation performance.* Same format as **Figure 2b** (except that reaction times are modeled in the left plot). (c) *Comparison of performance against control models.* Performance of the Drift diffusion model does not degrade substantially if a single threshold is used, thus justifying this simplification. Performance degrades if axis-aligned category vectors are used, supporting the assertion that responses of multiple VTC regions are used by subjects in deciding image category. (d) *Overall model architecture.* This schematic summarizes all components of our computational model (**Figures 2a, 5a, 6a**). Bottom-up visual information is encoded in the VTC fixation response (green box; Template model), fixation responses are routed to the IPS for evidence accumulation (purple box; Drift diffusion model), and then feedback from the IPS to VTC causes top-down modulation during the categorization and one-back tasks (yellow box; IPS-scaling model).

## MATERIALS AND METHODS

572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605

### Subjects

Eleven subjects participated in this study. Two subjects were excluded due to inability to identify VWFA in one subject and low signal-to-noise ratio in another subject, leaving a total of nine usable subjects (age range 25–32; six males, three females). All subjects were healthy right-handed monolingual native-English speakers, had normal or corrected-to-normal visual acuity, and were naive to the purposes of the experiment. Informed written consent was obtained from all subjects, and the experimental protocol was approved by the Washington University in St. Louis Institutional Review Board. Each subject participated in 1–3 scanning sessions, over the course of which anatomical data (T1-weighted high-resolution anatomical volume, diffusion-weighted MRI data) and functional data (retinotopic mapping, functional localizer, main experiment) were collected.

### Visual stimuli

Stimuli were presented using an NEC NP-V260X projector. The projected image was focused onto a backprojection screen and subjects viewed this screen via a mirror mounted on the RF coil. The projector operated at a resolution of  $1024 \times 768$  at 60 Hz, and the viewing distance was 340 cm. A Macintosh laptop controlled stimulus presentation using code based on Psychophysics Toolbox (Brainard, 1997; Pelli, 1997). Approximate gamma correction was performed by taking the square root of pixel intensity values before stimulus presentation. Behavioral responses were recorded using a button box.

The experiment consisted of 22 types of stimuli. All stimuli were small grayscale images (approximately  $2^\circ \times 2^\circ$ ) presented at fixation. Each stimulus type consisted of 10 distinct images (e.g. 10 different faces for a face stimulus), and a subset of these images were presented on each given trial.

*FACE.* This stimulus consisted of a face pictured from a frontal viewpoint. Ten distinct faces were prepared. Faces were masked using a circle with diameter  $2^\circ$ . The outer  $0.25^\circ$  of the mask was smoothly ramped using a cosine function.

*WORD.* This stimulus consisted of a 5-letter word. Ten distinct words were prepared. Letters were white on a gray background, generated using the Helvetica font, and occupied a rectangular region measuring  $3.15^\circ \times 1.05^\circ$ .

606  
607 *PHASE COHERENCE*. These stimuli consisted of the FACE and WORD stimuli prepared at four phase-  
608 coherence levels, 0%, 25%, 50%, and 75% (8 stimuli total). To achieve this, for each of the ten images  
609 from each stimulus type, the portion of the image within a fixed region (FACE:  $2^\circ \times 2^\circ$  square; WORD:  
610  $3.15^\circ \times 1.05^\circ$  rectangle) was extracted, and its phase spectrum was blended, to different degrees, with a  
611 randomly generated phase spectrum. For example, 25% coherence indicates that the phase of each Fourier  
612 component was set to a value that lies at 75% of the angular distance from the original phase to the phase  
613 in the randomly generated spectrum.

614  
615 *NOISE*. This stimulus is the same as the FACE stimulus at 0% phase coherence. For brevity, we refer to  
616 this stimulus as NOISE.

617  
618 *CONTRAST*. These stimuli consisted of the FACE, WORD, and NOISE stimuli prepared at three contrast  
619 levels (9 stimuli total). The contrasts of the original stimuli were taken to be 100%, and different contrast  
620 levels were achieved by scaling pixel intensity values towards the background value. Contrast levels of  
621 4%, 6%, and 10% were used for the FACE and NOISE stimuli, and contrast levels of 3%, 5%, and 8%  
622 were used for the WORD stimulus. This choice of contrast levels matches the average root-mean-square  
623 (RMS) contrast across stimulus types (e.g., the average RMS contrast for the FACE stimulus at 4%  
624 contrast is approximately equal to the average RMS contrast for the WORD stimulus at 3% contrast).  
625 Note that a contrast level of 0% was achieved by estimating responses to blank trials (see *GLM analysis*).

626  
627 *POLYGON*. This stimulus consisted of a string of three polygons (each chosen randomly from a set of  
628 polygons). Polygons were white, unfilled, on a gray background, and occupied a region similar in size to  
629 that of the WORD stimulus. Ten distinct strings were prepared.

630  
631 *CHECKERBOARD*. This stimulus consisted of alternating black and white square checks. Ten  
632 checkerboards were prepared by varying check size from  $0.03125^\circ$  to  $0.5^\circ$  using ten equally spaced steps  
633 on a logarithmic scale. The  $x$ - and  $y$ -positions of each checkerboard were set randomly. Checkerboards  
634 were masked using a circle with diameter  $2^\circ$ .

635  
636 *HOUSE*. This stimulus consisted of a house pictured from a frontal viewpoint. Ten distinct houses were  
637 prepared. Houses were masked using a  $2^\circ \times 2^\circ$  square. The outer  $0.25^\circ$  of the mask was smoothly ramped  
638 using a cosine function.

639

640 **Experimental design and tasks**

641  
642 Stimuli were presented in 4-s trials, one stimulus per trial. In a trial, four images from a given stimulus  
643 type (e.g. FACE, 10% contrast) were presented sequentially using an 800-ms ON, 200-ms OFF duty  
644 cycle. To generate the sequence of four images, we first randomly selected four distinct images out of the  
645 ten images associated with the stimulus type. Then, for certain trials (details below), we modified the  
646 sequence to include a repetition by randomly selecting one of the images (excluding the first) and  
647 replacing that image with the previous image. Throughout stimulus presentation, a small dot ( $0.12^\circ \times$   
648  $0.12^\circ$ ) was present at the center of the display. The dot switched to a new randomly selected color every  
649 600 ms using a set of six possible colors: magenta, red, yellow, green, cyan, and blue.

650  
651 In the experiment, two of the stimuli were duplicated (FACE and WORD), yielding a total of 24 stimulus  
652 conditions. Data corresponding to these duplicate stimuli are not used in this paper. Each run began and  
653 ended with a 16-s baseline period in which no stimuli were presented. During a run, each of the 24  
654 stimulus conditions was presented three times. Six blank trials (no stimulus) were also included. The  
655 order of stimulus and blank trials was random, subject to the constraints that blank trials could not occur  
656 first nor last, blank trials could not occur consecutively, and no stimulus condition could occur  
657 consecutively. During the baseline periods and blank trials, the small central dot was still present. A  
658 randomly selected two of the three trials associated with each stimulus condition were modified to include  
659 an image repetition (as described previously). Each run lasted 344 seconds (5.7 min).

660  
661 For each run, subjects were instructed to maintain fixation on the central dot while performing one of  
662 three tasks. In the *fixation task*, subjects were instructed to press a button whenever the central dot turned  
663 red, and were additionally reminded to not confuse the red and magenta colors. In the *categorization task*,  
664 subjects were instructed to report for each stimulus trial whether they perceived a word, a face, or neither  
665 (“other”). Responses were made using three different buttons, and subjects were reminded to make only  
666 one response for each 4-s trial. Note that it is possible that responses are made prior to the completion of  
667 the four images that comprise a trial. In the *one-back task*, subjects were instructed to press a button  
668 whenever an image was repeated twice in a row, and were informed that repetitions occurred only within  
669 stimulus trials and not across trials. Subjects were warned that although some stimuli are faint (low  
670 contrast), they should still try their best to perform the categorization and one-back tasks. Subjects were  
671 also informed that some trials are blank trials and that responses were not expected on these trials.  
672 Subjects were familiarized with the stimuli and tasks before the actual experiment was conducted.

673

674 Subjects performed each of the three tasks four times during the course of the experiment, yielding a total  
675 of 3 tasks  $\times$  4 runs = 12 runs. The physical stimulus sequence (including the temporal ordering of  
676 stimulus images and dot colors) was held constant across tasks. This was accomplished by generating  
677 four distinct stimulus sequences and cycling through the sequences and tasks. Specifically, the order of  
678 stimulus sequences was ABCD ABCD ABCD, where each letter corresponds to a distinct sequence, and  
679 the order of tasks was XYZ XYZ XYZ XYZ, where each letter corresponds to a distinct task. The order  
680 of tasks was counterbalanced across subjects. Each stimulus and task combination (e.g.  
681 CHECKERBOARD during one-back task) occurred a total of 3 trials  $\times$  4 runs = 12 times over the course  
682 of the experiment.

683

#### 684 **MRI data acquisition**

685

686 MRI data were collected at the Neuroimaging Laboratory at the Washington University in St. Louis  
687 School of Medicine using a modified 3T Siemens Skyra scanner and a 32-channel RF coil. For functional  
688 data, 28 oblique slices covering occipitotemporal cortex were defined: slice thickness 2.5 mm, slice gap 0  
689 mm, field-of-view 200 mm  $\times$  200 mm, phase-encode direction anterior-posterior. A T2\*-weighted, single-  
690 shot, gradient-echo EPI sequence was used: matrix size 80  $\times$  80, TR 2 s, TE 30 ms, flip angle 77°,  
691 nominal spatial resolution 2.5 mm  $\times$  2.5 mm  $\times$  2.5 mm. Fieldmaps were acquired for post-hoc correction  
692 of EPI spatial distortion. To achieve comprehensive coverage for localization of top-down effects, a  
693 whole-brain version of the protocol involving 58 slices and a multiband (Feinberg et al., 2010) factor of 2  
694 was used in three of the nine subjects. In addition to functional data, T1-weighted anatomical data  
695 (MPRAGE sequence, 0.8-mm resolution) and diffusion-weighted data (spin-echo EPI sequence, 2-mm  
696 resolution, 84 directions, *b*-values of 1,500 and 3,000) were acquired. The diffusion sequence was  
697 acquired twice, reversing the phase-encode direction, in order to compensate for spatial distortions.  
698 Diffusion data were not acquired for one subject due to time constraints.

699

#### 700 **Behavioral analysis**

701

702 Behavioral results for the categorization task are used in the present study. We analyzed both reaction  
703 times (RT) and category judgments. We defined RT as the time elapsed between the onset of the first of  
704 the four images in a given trial and the button press. Trials in which no buttons were pressed were  
705 ignored. For each subject, we summarized RTs by computing the median RT across trials for each  
706 stimulus. To obtain group-averaged RTs, we added a constant to each subject's RTs in order to match the  
707 mean RT to the grand mean across subjects and then computed the mean and standard error across

708 subjects (this normalization procedure compensates for additive offsets in RT across subjects). Category  
709 judgments were analyzed by calculating percentages of trials on which a given subject categorized a  
710 given stimulus into each of the three categories (word, face, other). Subjects were highly consistent in  
711 their judgments: for each stimulus, the most frequently reported category was the same across subjects  
712 and was reported more than 85% of the time. Category judgments obtained from the categorization task  
713 are used in the labeling and interpretation of experimental results (e.g. **Figures 1, 3**).

714

### 715 **Diffusion analysis**

716

717 Subject motion was corrected by co-registering each volume to the average of the non-diffusion-weighted  
718  $b = 0$  images. Gradient directions were adjusted to account for the co-registration. From pairs of volumes  
719 acquired with reversed phase-encode directions, the susceptibility-induced off-resonance field was  
720 estimated using a method similar to that described in Andersson et al., 2003 as implemented in FSL  
721 (Behrens et al., 2004). Eddy currents were corrected using FSL's *eddy* tool. The  $b = 3,000$  measurements  
722 were used to estimate fiber orientation distribution functions for each voxel using constrained spherical  
723 deconvolution as implemented in *mrtrix* (Tournier et al., 2007) (CSD,  $l_{max} = 4$ ), and fiber tracts were  
724 estimated using probabilistic tractography (500,000 fibers). For each subject, we identified the vertical  
725 occipital fasciculus (VOF) using a previously published algorithm (Yeatman et al., 2014), and then  
726 quantified the Euclidean distance from the VOF terminations to word- and face-selective regions in VTC  
727 and the task-related hotspot in the IPS.

728

### 729 **Pre-processing of anatomical and functional data**

730

731 The T1-weighted anatomical volume acquired for each subject was processed using FreeSurfer (Fischl,  
732 2012). The results were used to create a cortical surface reconstruction positioned halfway between the  
733 pial surface and the boundary between gray and white matter. We used the *fsaverage* surface from  
734 FreeSurfer to define anatomical ROIs (details below). These ROIs were transformed to native subject  
735 space by performing nearest-neighbor interpolation on the spherical surfaces created by FreeSurfer (these  
736 surfaces reflect folding-based alignment of individual subject surfaces to the *fsaverage* surface).

737

738 Functional data were pre-processed by performing slice time correction, fieldmap-based spatial  
739 undistortion, motion correction, and registration to the subject-native anatomical volume. The combined  
740 effects of distortion, motion, and registration were corrected using a single cubic interpolation of the slice

741 time corrected volumes. Interpolations were performed directly at the vertices of the subject's cortical  
742 surface, thereby avoiding unnecessary interpolation and improving spatial resolution (Kang et al., 2007).

743

#### 744 **GLM analysis**

745

746 The pre-processed fMRI data were analyzed using GLMdenoise (Kay et al., 2013a)  
747 (<http://kendrickkay.net/GLMdenoise/>), a data-driven denoising method that derives estimates of  
748 correlated noise from the data and incorporates these estimates as nuisance regressors in a general linear  
749 model (GLM) analysis of the data. For our experiment, we coded each stimulus and task combination as a  
750 separate condition and also included the blank trials, producing a total of  $(24 \text{ stimulus} + 1 \text{ blank}) \times 3 \text{ tasks}$   
751  $= 75$  conditions. The response to blank trials was interpreted as the response to a 0%-contrast stimulus.  
752 Estimates of BOLD response amplitudes (beta weights) were converted to units of percent BOLD signal  
753 change by dividing amplitudes by the mean signal intensity observed at each vertex. To obtain ROI  
754 responses, beta weights were averaged across the vertices composing each ROI. Error bars (68% CIs) on  
755 beta weights were obtained by bootstrapping runs.

756

757 Group-averaged beta weights were calculated using a procedure that compensates for large intrinsic  
758 variation in percent BOLD change across subjects. First, the beta weights obtained for each subject in a  
759 given ROI were normalized to be a unit-length vector (e.g.  $\tilde{\mathbf{b}}_i = \mathbf{b}_i / \|\mathbf{b}_i\|$  where  $\mathbf{b}_i$  indicates beta weights  
760 for the  $i$ th subject ( $1 \times n$ ),  $\|\cdot\|$  indicates  $L_2$ -norm, and  $\tilde{\mathbf{b}}_i$  indicates normalized beta weights for the  $i$ th  
761 subject). Next, normalized beta weights were averaged across subjects, using bootstrapping to obtain error  
762 bars (68% CIs). Finally, the resulting group-averaged beta weights were multiplied by a scalar such that  
763 the mean of the beta weights is equal to the mean of the original unnormalized beta weights obtained from  
764 all subjects. The motivation of this last step is to produce interpretable units of percent BOLD change  
765 instead of normalized units. Note that in some cases, beta weights are repeated for easier visualization  
766 (e.g., in **Figure 1**, NOISE at 100% contrast is the same data point as FACE at 0% phase coherence).  
767 Group-averaged beta weights were used in computational modeling.

768

#### 769 **Region-of-interest (ROI) definition**

770

771 Visual field maps were defined using the population receptive field (pRF) technique applied to retinotopic  
772 mapping data (Dumoulin and Wandell, 2008; Kay et al., 2013b). Subjects participated in 2–4 runs (300-s  
773 each) in which they viewed slowly moving apertures (bars, wedges, rings) filled with a colorful texture of

774 objects, faces, and words placed on an achromatic pink-noise background. The aperture and texture were  
775 updated at 5 Hz, and blank periods were included in the design (Dumoulin and Wandell, 2008). A semi-  
776 transparent fixation grid was superimposed on top of the stimuli (Schira et al., 2009). Stimuli occupied a  
777 circular region with diameter  $10^\circ$  and the viewing distance was 251 cm. A small semi-transparent central  
778 dot ( $0.15^\circ \times 0.15^\circ$ ) was present throughout the experiment and changed color every 1–5 s. Subjects were  
779 instructed to maintain fixation on the dot and to press a button whenever its color changed. The time-  
780 series data from this experiment were modeled using the Compressive Spatial Summation model (Kay et  
781 al., 2013b) as implemented in analyzePRF (<http://kendrickkay.net/analyzePRF/>). Angle and eccentricity  
782 estimates provided by the model were then visualized on cortical surface reconstructions and used to  
783 define V1, V2, V3, and hV4 (Brewer et al., 2005). Due to the limited amount of pRF data acquired, there  
784 was insufficient signal-to-noise ratio to define visual field maps in parietal cortex.

785  
786 Category-selective regions FFA and VWFA were defined using functional localizers (Weiner and Grill-  
787 Spector, 2010; 2011). Subjects participated in 2 runs (336-s each) in which they viewed blocks of words,  
788 faces, abstract objects, and noise patterns. Each block lasted 16 s and consisted of 16 images presented at  
789 a rate of 1 Hz. The images differed from those in the main experiment. In each run, the four stimulus  
790 types were presented four times each in pseudorandom order, with occasional 16-s blank periods. A semi-  
791 transparent fixation grid was superimposed on top of the stimuli (Schira et al., 2009). Stimuli occupied a  
792  $4^\circ \times 4^\circ$  square region, with the words, faces, and objects occupying the central  $3^\circ \times 3^\circ$  of this region. The  
793 viewing distance was 340 cm. Subjects were instructed to maintain central fixation and to press a button  
794 when the same image is presented twice in a row. The time-series data from this experiment were  
795 analyzed using a GLM to estimate the amplitude of the BOLD response to the four stimulus categories.

796  
797 To define FFA and VWFA, responses to the four stimulus categories were visualized on cortical surface  
798 reconstructions. FFA and VWFA were defined based on stimulus selectivity, anatomical location, and  
799 topological relationship to retinotopic areas (Weiner and Grill-Spector, 2010; Yeatman et al., 2013;  
800 Weiner et al., 2014). We defined FFA as face-selective cortex (responses to faces greater than the average  
801 response to the other three categories) located on the fusiform gyrus. We included in the definition both  
802 the posterior fusiform gyrus (pFus-faces/FFA-1) and middle fusiform gyrus (mFus-faces/FFA-2)  
803 subdivisions of FFA (Weiner et al., 2014). We defined VWFA as word-selective cortex (responses to  
804 words greater than the average response to the other three categories) located in and around the left  
805 occipitotemporal sulcus. In some subjects, multiple word-selective patches were found, and all of these  
806 patches were included in the definition of VWFA.

807

808 Anatomically-defined ROIs were also created (see **Figures 4a** and **4b**). Based on curvature values on the  
809 *fsaverage* surface, we created an anatomical mask of the IPS by selecting the posterior segment of the  
810 intraparietal sulcus (Pitzalis et al., 2012). Using the atlas of visual topographic organization provided by  
811 Wang et al. (Wang et al., 2014), we estimate that this IPS mask overlaps V3A, V3B, IPS-0, IPS-1, and  
812 IPS-2. The locations of IPS-0/1/2/3 from the atlas are shown in **Figure 4** and **Figure 4-figure**  
813 **supplement 1**. We also created an anatomical mask of VTC by computing the union of the *fusiform* and  
814 *inferiortemporal* parcels provided by the FreeSurfer Desikan-Killiany atlas (Desikan et al., 2006) and  
815 trimming the anterior extent of the result to include only visually responsive cortex. The VTC mask  
816 includes both FFA and VWFA as well as surrounding cortex.

817

818 In our data, we find that word-selective visual cortex in some subjects is confined to the left hemisphere,  
819 consistent with previous studies (Yeatman et al., 2013). Therefore, to ease interpretation, we restricted our  
820 analysis to VWFA, FFA, VTC, and IPS taken from the left hemisphere. In addition, we restricted the  
821 definition of V1, V2, V3, hV4, VTC, and IPS to include only vertices exhibiting response amplitudes in  
822 the main experiment that are positive on average. This procedure excludes voxels with peripheral  
823 receptive fields which typically exhibit negative BOLD responses to centrally presented stimuli.

824

### 825 **Task-based functional connectivity**

826

827 To identify the cortical region that generates top-down effects in VWFA and FFA, we performed a simple  
828 connectivity analysis. First, we averaged BOLD responses across our VTC mask, given that top-down  
829 effects appear broadly across VTC. Next, we identified the component of the VTC response that is of no  
830 interest, specifically, the bottom-up stimulus-driven response. Our estimate of this component is given by  
831 our measurement of VTC responses during the fixation task (22 stimuli + 1 blank = 23 values). We then  
832 subtracted the bottom-up component from the VTC response measured during the categorization task (23  
833 values) and one-back task (23 values). This produced a set of residuals (46 values) that reflect the top-  
834 down effect in VTC. Finally, we correlated the residuals with the responses of each cortical location in  
835 our dataset during the categorization and one-back tasks (46 values). The cortical location that best  
836 correlates with the residuals is interpreted as a candidate region that supplies top-down modulation to  
837 VTC.

838

839 Results were visualized by averaging correlation values across subjects based on the *fsaverage* cortical  
840 alignment and plotting results on the *fsaverage* surface. Results from the three subjects for which whole-  
841 brain fMRI data were acquired are shown in **Figures 4a** and **4b**. Results from the remaining six subjects

842 with limited fMRI coverage are provided in **Figure 4–figure supplement 1**. Note that the correlation-  
843 based analysis we have used is most suitable for connectivity effects that are additive in nature (e.g., IPS  
844 providing additive enhancement to VTC). However, the modulation is more accurately characterized as a  
845 multiplicative, or scaling, effect (see **Figures 3b** and **3c**). The advantage of correlation is that it is robust  
846 to noise and computationally efficient; we perform a more precise evaluation of different top-down  
847 mechanisms in the computational modeling section below.

848

849 There are three important differences between the connectivity analysis described here and conventional  
850 correlation-based resting-state functional connectivity (RSFC) (Buckner et al., 2013) and the  
851 psychophysiological interactions (PPI) technique (O'Reilly et al., 2012). One is that our connectivity is  
852 performed on data that have explicit manipulation of stimulus and task (unlike RSFC). Another is that we  
853 analyze the data explicitly in terms of information-processing operations performed by the brain (unlike  
854 PPI). In other words, functional connectivity is characterized, not as correlated signal fluctuations, but as  
855 a direct consequence of information-processing operations. A third difference is that our connectivity is  
856 performed on beta weights that pool across trials (Rissman et al., 2004), as opposed to raw BOLD time-  
857 series. This concentrates the analysis on brain responses that are reliably driven by the stimulus and task,  
858 and de-emphasizes trial-to-trial fluctuations in cognitive performance (Donner et al., 2013).

859

## 860 **Computational modeling**

861

862 We developed a computational model to account for BOLD responses measured in VTC and IPS. The  
863 model is composed of three components, each of which addresses a different aspect of the data (**Figure**  
864 **6d**). The first component (*Template model*) specifies how a given stimulus drives bottom-up VTC  
865 responses as measured during the fixation task; the second component (*IPS-scaling model*) specifies how  
866 top-down modulation from the IPS during the categorization and one-back tasks affects VTC responses;  
867 and the third component (*Drift-diffusion model*) specifies how accumulation of evidence from VTC  
868 predicts reaction times and IPS responses during the categorization task. Note that although the three  
869 model components could be yoked together (e.g., the output from the Template model could serve as the  
870 input to the Drift-diffusion model), in our model implementations, we adopt the approach of isolating  
871 each model component so that the quality of each component can be assessed independently of the others.

872

873 For all three model components, computational modeling was performed using nonlinear least-squares  
874 optimization (MATLAB Optimization Toolbox). Leave-one-stimulus-out cross-validation was used to  
875 assess model accuracy (thus, we assess the ability of models to generalize to stimuli that the models have

876 not been trained on). Note that the use of cross-validation enables fair comparison of models that have  
877 different levels of flexibility (or, informally, different numbers of free parameters). This is because  
878 models that are overly complex will tend to fit noise in the training data and thereby generalize poorly to  
879 the testing data (Hastie et al., 2001).

880  
881 Accuracy was quantified as the percentage of variance explained ( $R^2$ ) between cross-validated predictions  
882 of the data (aggregated across cross-validation iterations) and the actual data. In the case of beta weights,  
883 variance was computed relative to 0% BOLD signal change (Kay et al., 2013b). In certain cases, accuracy  
884 is reported using Pearson's correlation ( $r$ ); this metric assesses performance relative to the mean. To  
885 assess reliability of cross-validation results, model fitting and cross-validation were repeated for each  
886 bootstrap of the group-averaged data (resampling subjects with replacement). For benchmarks on cross-  
887 validation performance, we calculated noise ceilings using Monte Carlo simulations (Kay et al., 2013b)  
888 and quantified the performance of a flat-response model that predicts the same response level for each  
889 data point.

890

### 891 **Template model**

892

893 *Basic model description.* The *Template model* specifies the stimulus properties that drive bottom-up  
894 responses in VTC. The model accepts as input a grayscale image and produces as output the predicted  
895 response in VWFA and FFA during the fixation task. In brief, the model processes the image using a set  
896 of V1-like Gabor filters and then computes a normalized dot product between filter outputs and a category  
897 template. The category template can be viewed as capturing the prototypical image statistics of a word  
898 (VWFA) or face (FFA). The Template model makes no claim as to how the brain might develop category  
899 templates; they might be genetically hard-wired (Kanwisher, 2010) or arise from experience with the  
900 environment (Gauthier et al., 1999). The central claim is that the bottom-up information computed by  
901 VTC is, at least to a first approximation, the output of a template operation applied to the stimulus.

902

903 The Template model is related to our previously developed Second-order contrast (SOC) model (Kay et  
904 al., 2013c). Similar to the SOC model, the Template model has a cascade architecture involving two  
905 stages of filtering, rectification, and normalization. The first stage of the Template model is taken directly  
906 from the SOC model, and the properties of this stage (e.g. filter design) were not tweaked to fit the data.  
907 The main difference between the Template and SOC models is that the Template model incorporates a  
908 specific second-stage filter (the template), whereas the SOC model uses a variance-like operation in the  
909 second stage that captures generic sensitivity to second-order contrast. Whether the Template model

910 captures certain response properties, such as invariance to font in VWFA (Dehaene and Cohen, 2011) or  
 911 coarse luminance-contrast selectivity in FFA (Ohayon et al., 2012), is an empirical question that can only  
 912 be resolved through quantitative evaluation on experimental data. For example, our measurements  
 913 indicate that VWFA responds strongly to polygons (**Figure 1d**); the Template model already accounts for  
 914 this effect (**Figure 2b**).

915

916 *Model details.* The first stage of the Template model involves computing a V1-like representation of the  
 917 image. The image is first resized to 250 pixels  $\times$  250 pixels, and luminance values are mapped to the  
 918 range  $[-0.5, 0.5]$ , which has the effect of mapping the gray background to 0. The model then calculates V1  
 919 energy in the same way as the SOC model (Kay et al., 2013c). Specifically, the image is projected onto a  
 920 set of isotropic Gabor filters occurring at 8 orientations, 2 quadrature phases, and a range of positions (63  
 921  $x$ -positions  $\times$  63  $y$ -positions). Filters are constructed at a single scale with a peak spatial frequency tuning  
 922 of 4 cycles per degree (see **Figure 2d**) and a spatial frequency bandwidth of 1 octave (full-width at half-  
 923 maximum of the amplitude spectrum). Filters are scaled such that filter responses to full-contrast optimal  
 924 sinusoidal gratings are equal to one. Outputs of quadrature-phase filters are squared, summed, and square-  
 925 rooted, analogous to the complex-cell energy model (Adelson and Bergen, 1985).

926

927 After computing V1 energy, the model applies divisive normalization (Heeger, 1992), again analogous to  
 928 the SOC model. The output of each filter is divided by the average output across filter orientations at the  
 929 same position:

$$930 \quad ncc_{pos,or} = \frac{(cc_{pos,or})^r}{s^r + \left( \frac{\sum_{or} cc_{pos,or}}{numor} \right)^r} \quad (1)$$

931 where  $ncc_{pos,or}$  is the normalized filter output at a given position and orientation,  $cc_{pos,or}$  is the filter output  
 932 at a given position and orientation,  $numor$  is the total number of orientations, and  $r$  and  $s$  are parameters  
 933 that control the strength of the normalization. For simplicity and to reduce the potential for overfitting, we  
 934 do not fit  $r$  and  $s$  but simply use  $r = 1$  and  $s = 0.5$ , values determined from our previous study (Kay et al.,  
 935 2013c).

936

937 At this point in the model, the representation of the image is a 3D matrix of dimensions 63  $x$ -positions  $\times$   
 938 63  $y$ -positions  $\times$  8 orientations. To visualize this representation, a hue-saturation-value image is used (see  
 939 **Figure 2a**). For each position, a set of 8 vectors is constructed with vector angles corresponding to filter

940 orientation and vector lengths corresponding to normalized filter output. These vectors are averaged and  
941 an image pixel is used to summarize the result. Specifically, the hue of a pixel indicates the angle of the  
942 vector average and the value of the pixel indicates the length of the vector average.

943

944 The second stage of the Template model involves taking the V1-like representation of the image and  
945 comparing it to a category template to generate the predicted response. Specifically, the response is  
946 computed as

$$947 \quad RESP = a \times \frac{|S \cdot T - b\bar{S}|^+}{c + \bar{S}} \quad (2)$$

948 where  $RESP$  is the predicted response,  $S$  is the 3D matrix with the V1-like representation of the image,  $T$   
949 is the category template,  $\bar{S}$  is the average of the elements in  $S$ ,  $| \cdot |^+$  indicates positive half-wave  
950 rectification, and  $a$ ,  $b$ , and  $c$  are free parameters (3 free parameters).

951

952 There are three basic steps in **Equation 2**. The first step is a filtering operation, accomplished by  
953 computing the dot product between the stimulus and the template ( $S \cdot T$ ). Intuitively, this operation  
954 quantifies the similarity between the stimulus and the template. The second step is subtraction of average  
955 stimulus energy ( $-b\bar{S}$ ) with a free parameter controlling the strength of the subtractive normalization.  
956 This subtraction can be interpreted as penalizing non-specific energy in the stimulus, thereby inducing  
957 preference for stimulus energy that conforms to the category template. (An alternative interpretation is  
958 that the subtraction provides flexibility with respect to the overall mean of the template.) To ease  
959 interpretation and ensure that negative responses are not obtained, the result of the subtraction is  
960 positively rectified ( $| \cdot |^+$ ). The third step is division by average stimulus energy ( $(c + \bar{S})$ ) with a free  
961 parameter controlling the strength of the divisive normalization. This division penalizes non-specific  
962 energy in the stimulus, similar to subtractive normalization, but induces a different response geometry  
963 (Zetsche et al., 1999). In summary, **Equation 2** computes a dot product between the stimulus and the  
964 template that is normalized subtractively and divisively by the average stimulus energy.

965

966 Where does the category template in **Equation 2** come from? Given that we do not have sufficient  
967 sampling of stimuli to directly estimate templates from the data (but see **Figure 2–figure supplement 1**),  
968 we adopted the simple strategy of constructing templates from our stimulus set. Specifically, we took the  
969 WORD and FACE stimuli at 100% contrast and used the first stage of the Template model to compute a  
970 V1-like representation of these stimuli. This produced for each category, ten points in a  $63 \times 63 \times 8 =$   
971 31,752-dimensional space. We then computed the centroid of the ten points, producing a category

972 template (example shown in **Figure 2a**). Because the category template is constructed from the same  
973 stimuli used in our experiment, it is guaranteed that the Template model predicts large responses to the  
974 preferred category (e.g., using a category template constructed from the face stimuli guarantees that the  
975 face stimuli produce large responses from the model). However, there is no guarantee that the model will  
976 accurately account for responses to the other stimuli used in our experiment.

977  
978 *Model fitting.* The Template model was fit to the fixation responses of VWFA and FFA. Model outputs  
979 were calculated for all ten images associated with a given stimulus type and then averaged to obtain the  
980 final model output for that stimulus type. To aid model fitting, the  $S \cdot T$  and  $\bar{S}$  quantities were pre-  
981 computed and pre-conditioned by dividing each quantity by the mean of that quantity across stimuli. After  
982 pre-conditioning, a variety of initial seeds for  $b$  and  $c$  were evaluated in order to avoid local minima.  
983 Specifically, we performed optimization starting from initial seeds corresponding to every combination of  
984  $b$  and  $c$ , where  $b$  is chosen from  $\{0.5, 1, 1.5, 2, 3, 5\}$  and  $c$  is chosen from  $\{.01, .05, .1, .5, 1, 5, 10\}$ .

985  
986 *Alternative models.* (1) The *Category model* predicts a fixed response level for stimuli from the preferred  
987 stimulus category (word for VWFA, face for FFA) and a different response level for all other stimuli (2  
988 free parameters, one for each response level). Category judgments provided by the subjects were used to  
989 determine category membership; for example, words and faces at 0% and 25% phase coherence were  
990 reported by subjects as ‘other’, and are hence not considered to be words and faces by the Category  
991 model. (2–3) We evaluated simplified versions of the second-stage normalization used in the Template  
992 model. One version, *Template model (only subtractive normalization)*, omits the divisive normalization  
993 and thus characterizes responses as a simple linear function of V1-like normalized filter outputs (2 free  
994 parameters,  $a$  and  $b$ ), whereas the other version, *Template model (only divisive normalization)*, omits the  
995 subtractive normalization (2 free parameters,  $a$  and  $c$ ). (4) In *Template model (omit first stage)*, the first  
996 stage of the model is omitted and the template operation is performed on a pixel representation of the  
997 image, i.e.,  $S$  refers to the original image instead of the V1-like representation of the image (3 free  
998 parameters,  $a$ ,  $b$ , and  $c$ ). (5–7) We evaluated the effect of using different templates in the Template model  
999 (each model has 3 free parameters,  $a$ ,  $b$ , and  $c$ ). *Template model (non-selective template)* uses a template  
1000 consisting of all ones. *Template model (mixed template)* uses a template generated by unit-length  
1001 normalizing both the word and face templates and then averaging the templates together. *Template model*  
1002 *(random template)* uses a template generated by drawing uniform random values from the range  $[0,1]$ .

1003  
1004 **IPS-scaling model**

1005

1006 *Basic model description.* The *IPS-scaling model* predicts top-down modulation of VTC by taking into  
1007 account measurements of IPS activity. The model accepts as input the response in VTC (either VWFA or  
1008 FFA) during the fixation task and the response in IPS during the stimulus-directed tasks (categorization,  
1009 one-back), and produces as output the predicted response in VTC during the stimulus-directed tasks.  
1010 Intuitively, the model answers the question: how much is the bottom-up response in VTC enhanced by the  
1011 IPS when the subject performs a task on the stimulus? The model can be viewed as a formal  
1012 implementation of the concept of stimulus-specific scaling (schematized in **Figure 3b, lower right**).  
1013 Similar ideas regarding top-down scaling induced by the IPS can be found in previous work (Kayser et  
1014 al., 2010b).

1015

1016 *Model details.* The IPS-scaling model multiplies the bottom-up response in VTC measured during the  
1017 fixation task by a scaled version of the IPS response observed during a stimulus-directed task:

$$1018 \quad VTC_{task} = VTC_{bottom} \times (a \cdot IPS_{task} + b) \quad (3)$$

1019 where  $VTC_{task}$  is the predicted response in VTC during the stimulus-directed task,  $VTC_{bottom}$  is the bottom-  
1020 up response in VTC,  $IPS_{task}$  is the response in IPS during the stimulus-directed task, and  $a$  and  $b$  are  
1021 parameters that allow a scale and offset to be applied to the IPS response (2 free parameters). The final  
1022 scaling factor that is applied to  $VTC_{bottom}$  is shown in **Figure 5b**. The measurements of IPS activity used in  
1023 the model are extracted using a broad anatomical mask of the IPS (see *Region-of-interest (ROI)*  
1024 *definition*).

1025

1026 *Model fitting.* The IPS-scaling model was fit to the fixation, categorization, and one-back responses  
1027 observed in VWFA and FFA. Leave-one-out cross-validation was performed by systematically leaving  
1028 out each of the categorization and one-back responses. Since measurement noise is present in the fixation  
1029 responses, treating the fixation responses as exact estimates of bottom-up responses would result in  
1030 suboptimal model performance (especially in the case of bottom-up responses that are near zero). We  
1031 therefore devised a procedure that allows flexibility in estimating bottom-up responses (see light lines in  
1032 **Figure 5b**). In the procedure, a separate parameter is used to model the bottom-up response associated  
1033 with each stimulus. During model fitting, these bottom-up parameters are initially set to be equal to the  
1034 measured fixation responses, parameters of the model excluding the bottom-up parameters are optimized,  
1035 and then all parameters are optimized simultaneously. This procedure was also used for the alternative  
1036 models described below. Note that the IPS-scaling model uses flexible parameters to accommodate  
1037 bottom-up stimulus selectivity and does not attempt to characterize the image-processing computations  
1038 that underlie bottom-up responses (such computations are in the purview of the Template model).

1039

1040 *Alternative models.* (1) The *Task-invariant model* posits that top-down modulation does not occur and  
1041 that a fixed set of responses can characterize all three tasks (0 free parameters). (2–5) We evaluated  
1042 several phenomenological models for purposes of comparison. The *Additive model* (schematized in  
1043 **Figure 3b, upper left**) predicts responses during stimulus-directed tasks by adding a constant to bottom-  
1044 up responses (1 free parameter). The *Scaling model* (schematized in **Figure 3b, lower left**) predicts  
1045 responses during stimulus-directed tasks by multiplying bottom-up responses by a constant (1 free  
1046 parameter). The *Additive model (task-specific)* and *Scaling model (task-specific)* are identical to the  
1047 previous two models, except that separate constants are used for the categorization and one-back tasks (2  
1048 free parameters). (6) The *Area-specific enhancement model* (schematized in **Figure 3b, upper right**) is  
1049 identical to the *Scaling model (task-specific)* except that scaling is applied only to the stimuli preferred by  
1050 a given area, i.e. words in VWFA and faces in FFA (2 free parameters). (7) The *IPS-additive model*  
1051 predicts responses during stimulus-directed tasks by adding a scaled version of the IPS response to  
1052 bottom-up responses in VTC (2 free parameters). (8–9) To assess the specificity of the IPS enhancement,  
1053 we evaluated variants of the IPS-scaling model. In the *IPS-scaling (shuffle)* model, IPS responses are  
1054 shuffled across stimuli and tasks (restricted to the stimulus-directed tasks) before being used in the model  
1055 (2 free parameters). In the *IPS-scaling (shuffle within task)* model, IPS responses are shuffled across  
1056 stimuli but not across tasks before being used in the model (2 free parameters).

1057

### 1058 **Drift diffusion model**

1059

1060 *Basic model description.* The *Drift diffusion model* specifies the decision-making operations that underlie  
1061 performance of the categorization task, and is based upon past research on perceptual decision-making  
1062 (Shadlen and Newsome, 2001; Heekeren et al., 2004; Gold and Shadlen, 2007; Kayser et al., 2010a). The  
1063 model accepts as input fixation responses in VTC and produces as output predicted reaction times and IPS  
1064 responses for the categorization task. The basic idea is that VTC responses provide evidence regarding  
1065 which stimulus category has been presented to the subject, and this evidence is accumulated over time by  
1066 the IPS in order to make a final decision regarding stimulus category.

1067

1068 *Model details.* First, we collect fixation responses in hV4, VWFA, and FFA and divide each set of  
1069 responses by their mean. This normalization ensures that different ROIs have similar units. Then, for each  
1070 stimulus category (word, face, other), we compute the centroid of the fixation responses associated with  
1071 that category, interpret this centroid as a vector, and normalize the vector to unit length. This procedure  
1072 generates category vectors, defined in a three-dimensional neural space, that point in the directions of the  
1073 "arms" of the manifold of the VTC representation (see **Figure 3e**).

1074

1075 Next, we take the VTC fixation response for a given stimulus and project this response onto the category  
1076 vector associated with that stimulus. The working hypothesis is that this operation is performed by  
1077 neurons in IPS and that the magnitude of the projection indicates the strength of evidence for that specific  
1078 category. For example, there might be an IPS neuron that receives information from VTC and responds  
1079 strongly when the VTC response is consistent with the category vector corresponding to a word.

1080

1081 In accordance with drift diffusion models, we posit that evidence is accumulated until a threshold is  
1082 reached, at which point the decision is made. This generates a prediction of the reaction time required to  
1083 perform the categorization task on the stimulus:

$$1084 \quad RT = c + \frac{b}{R \cdot V} \quad (4)$$

1085 where  $RT$  is the predicted reaction time,  $R$  is the VTC fixation response,  $V$  is the category vector  
1086 associated with the stimulus,  $b$  is a parameter that controls the threshold, and  $c$  is a parameter that  
1087 compensates for non-decision time (e.g. motor response) (2 free parameters).  $R \cdot V$  is interpreted as a drift  
1088 rate, and  $b/(R \cdot V)$  is the time required to reach the threshold (see **Figure 6a**). Note that our instantiation of  
1089 the drift diffusion model is relatively simple, as it is non-stochastic and does not characterize trial-to-trial  
1090 variability. Thus, it can be viewed as a simplified version of the classic drift diffusion model (Ratcliff,  
1091 1978). Also, being non-stochastic, our model bears similarity to the linear ballistic accumulator model  
1092 (Brown and Heathcote, 2008), a model that also uses the idea of evidence accumulation (see Donkin et  
1093 al., 2011 for discussion of these different models).

1094

1095 Given that neuronal responses in parietal cortex reflect the duration of the decision-making process  
1096 (Shadlen and Newsome, 2001), we can use  $RT$  to predict IPS activity. A detailed model relating  $RT$  to  
1097 BOLD measurements of IPS activity requires precise characterization of neural dynamics during  
1098 decision-making and IPS subdivisions that might represent evidence accumulation for different stimulus  
1099 categories. For the purposes of this study, we use a simple model that posits a monotonically increasing  
1100 relationship between  $RT$  and the IPS response:

$$1101 \quad IPS = a \times \tanh(b \cdot RT + c) + d \quad (5)$$

1102 where  $IPS$  is the predicted IPS response,  $RT$  is the observed reaction time for a given stimulus,  $\tanh$  is the  
1103 hyperbolic tangent function intended as a generic sigmoidal nonlinearity, and  $a$ ,  $b$ ,  $c$ , and  $d$  are free  
1104 parameters (4 free parameters).

1105

1106 *Alternative models.* (1) The *Drift diffusion model (separate thresholds)* uses a separate threshold  
1107 parameter for each stimulus category (4 free parameters, one for non-decision time and three for  
1108 thresholds). This allows us to assess the validity of having a single threshold parameter in the model. (2)  
1109 The *Drift diffusion model (axis-aligned category vectors)* uses category vectors that are aligned with the  
1110 axes of the multi-dimensional neural space (4 free parameters, similar to the previous model). For  
1111 example, in this model, the word category vector is a vector that is one along the VWFA axis and zero  
1112 along the hV4 and FFA axes. This model tests the idea that evidence for words and faces is contributed  
1113 only by the VTC regions selective for those categories.

1114

1115 **Additional wide-range-of-stimuli dataset**

1116

1117 *Experimental design.* To assess the generalization performance of the Template model, we collected an  
1118 additional dataset involving a wider range of stimuli. This dataset was collected from one subject (an  
1119 author; male; age 34). Informed written consent was obtained, and the protocol was approved by the  
1120 University of Minnesota Institutional Review Board. The experiment was similar in design to the main  
1121 experiment. Stimuli included 22 images from the main experiment (one image from each of the 22  
1122 stimulus types), 92 images from a previous study investigating object representation in ventral temporal  
1123 cortex (Kriegeskorte et al., 2008), and 19 other images not used in this paper. As in the main experiment,  
1124 images were approximately  $2^\circ \times 2^\circ$  in size. In each 4-s trial, a single image was flashed using an 800-ms  
1125 ON, 200-ms OFF duty cycle. During a run, each image was presented in one trial, and 11 blank trials  
1126 were also included. Each run lasted 608 seconds (10.1 min), and a total of 10 runs were collected. During  
1127 stimulus presentation, the subject performed a variant of the fixation task. A small dot ( $0.1^\circ \times 0.1^\circ$ ) was  
1128 present at the center of the display and switched to one of five shades of red (ranging from (40,0,0) to  
1129 (255,0,0) in five equally spaced increments) every 1200 ms (repetitions allowed). The subject was  
1130 instructed to press a button whenever the luminance of the central dot increased and a different button  
1131 whenever the luminance decreased.

1132

1133 *MRI data acquisition.* MRI data were collected at the Center for Magnetic Resonance Research at the  
1134 University of Minnesota using a 7T Siemens Magnetom scanner and a custom 4-channel-transmit, 32-  
1135 channel-receive RF head coil. Stimuli were presented using a Cambridge Research Systems BOLDscreen  
1136 32 LCD monitor (resolution  $1920 \times 1080$  at 120 Hz; viewing distance 189.5 cm). Functional data were  
1137 acquired using 84 oblique slices covering occipitotemporal cortex: slice thickness 0.8 mm, slice gap 0  
1138 mm, field-of-view  $160 \text{ mm (FE)} \times 129.6 \text{ mm (PE)}$ , phase-encode direction inferior-superior. A T2\*-  
1139 weighted, single-shot, gradient-echo EPI sequence was used: matrix size  $200 \times 162$ , TR 2.2 s, TE 22.4

1140 ms, flip angle 80°, phase partial Fourier 6/8, in-plane acceleration factor (iPAT) 3, slice acceleration  
1141 factor (multiband (Moeller et al., 2010)) 2, nominal spatial resolution 0.8 mm × 0.8 mm × 0.8 mm.

1142  
1143 *Data analysis.* After pre-processing, the functional data were averaged across the thickness of gray matter  
1144 and then analyzed using GLMdenoise (as in the main experiment). Beta weights extracted from FFA and  
1145 VWFA were then modeled using several variants of the Template model. Model accuracy was quantified  
1146 using 20-fold cross-validation (random subsets of the stimuli). (1) The *Template model (original)* is the  
1147 same model used in the main experiment (3 free parameters). Importantly, the category template in the  
1148 model is fixed and not adjusted to the new dataset. (2) The *Template model (half-max average)* is a  
1149 simple extension of the Template model in which the category template is estimated as follows: on each  
1150 cross-validation iteration, using only the training set, identify responses that are at least half of the  
1151 maximum response and then compute the centroid (in the V1-like representation) of the stimuli  
1152 corresponding to these responses (3 free parameters plus nonparametric fitting of the template). Note that  
1153 this procedure is cross-validated in the sense that the category template is fit only to the training set;  
1154 whether the estimated category template generalizes to novel stimuli is an empirical question that is  
1155 assessed through cross-validation. (3) The *Template model (half-max cluster)* further extends the model to  
1156 accommodate multiple category templates. The logic is that just as V1 models use filters at multiple  
1157 orientations and spatial scales to characterize the overall V1 response, we might conceptualize FFA and  
1158 VWFA as containing multiple templates tuned to different types of stimuli (e.g. different face viewpoints,  
1159 different fonts). First, we identify responses that are at least half of the maximum response from the  
1160 training set. We then perform *k*-means clustering (in the V1-like representation) on the stimuli  
1161 corresponding to these responses. We use a cosine metric to quantify distance, and we select the solution  
1162 that minimizes cluster assignment error across 100 random initializations of centroid positions. The  
1163 obtained centroids are unit-length normalized and then used as category templates in the Template model.  
1164 For simplicity and to avoid overfitting, we compute the predicted response as a simple sum across the  
1165 independent responses of different category templates and use the same parameter values (*a*, *b*, *c* in  
1166 Equation 2) for different templates. We systematically vary the number of clusters from 1 through 8, and  
1167 select the number that maximizes cross-validation performance (4 free parameters—three for *a*, *b*, and *c*,  
1168 one for the number of clusters—plus nonparametric fitting of templates).

1169  
1170 **Code availability**

1171  
1172 Software code implementing the model proposed in this paper is available at  
1173 <http://cvnlab.net/vtcipsmodel/>.

## REFERENCES

- 1174
- 1175 Adelson EH, Bergen JR (1985) Spatiotemporal energy models for the perception of motion. *Journal of the*  
1176 *Optical Society of America* 2:284–299.
- 1177 Andersson JLR, Skare S, Ashburner J (2003) How to correct susceptibility distortions in spin-echo echo-  
1178 planar images: application to diffusion tensor imaging. *NeuroImage* 20:870–888.
- 1179 Avidan G, Harel M, Hendler T, Ben-Bashat D, Zohary E, Malach R (2002) Contrast sensitivity in human  
1180 visual areas and its relationship to object recognition. *Journal of neurophysiology* 87:3102–3116.
- 1181 Baldauf D, Desimone R (2014) Neural mechanisms of object-based attention. *Science* 344:424–427.
- 1182 Behrens T, Beckmann C, Smith SM, Jenkinson M, Woolrich MW, Johansen-Berg H, Bannister PR, De  
1183 Luca M, Drobnjak I, Flitney DE, Niazy RK, Saunders J, Vickers J, Zhang Y, de Stefano N, Brady  
1184 JM, Matthews PM (2004) Advances in functional and structural MR image analysis and  
1185 implementation as FSL. *NeuroImage* 23 Suppl 1:S208–S219.
- 1186 Bejjanki VR, Beck JM, Lu Z-L, Pouget A (2011) Perceptual learning as improved probabilistic inference  
1187 in early sensory areas. *Nature neuroscience* 14:642–648.
- 1188 Boynton GM (2009) A framework for describing the effects of attention on visual responses. *Vision*  
1189 *research* 49:1129–1143.
- 1190 Brainard DH (1997) The Psychophysics Toolbox. *Spat Vis* 10:433–436.
- 1191 Brewer AA, Liu J, Wade AR, Wandell B (2005) Visual field maps and stimulus selectivity in human  
1192 ventral occipital cortex. *Nature neuroscience* 8:1102–1109.
- 1193 Brouwer GJ, Heeger DJ (2013) Categorical clustering of the neural representation of color. *J Neurosci*  
1194 33:15454–15465.
- 1195 Brown SD, Heathcote A (2008) The simplest complete model of choice response time: linear ballistic  
1196 accumulation. *Cogn Psychol* 57:153–178.
- 1197 Buckner RL, Krienen FM, Yeo BTT (2013) Opportunities and limitations of intrinsic functional  
1198 connectivity MRI. *Nature neuroscience* 16:832–837.
- 1199 Buracas GT, Boynton GM (2007) The effect of spatial attention on contrast response functions in human  
1200 visual cortex. *J Neurosci* 27:93–97.
- 1201 Carandini M, Demb JB, Mante V, Tolhurst DJ, Dan Y, Olshausen BA, Gallant JL, Rust NC (2005) Do we  
1202 know what the early visual system does? *J Neurosci* 25:10577–10597.
- 1203 Cohen L, Dehaene S, Naccache L, Lehéricy S, Dehaene-Lambertz G, Hénaff MA, Michel F (2000) The  
1204 visual word form area: spatial and temporal characterization of an initial stage of reading in normal  
1205 subjects and posterior split-brain patients. *Brain* 123 ( Pt 2):291–307.
- 1206 Cohen L, Lehéricy S, Chochon F, Lemer C, Rivaud S, Dehaene S (2002) Language-specific tuning of  
1207 visual cortex? Functional properties of the Visual Word Form Area. *Brain* 125:1054–1069.
- 1208 Corbetta M, Shulman GL (2002) Control of goal-directed and stimulus-driven attention in the brain.

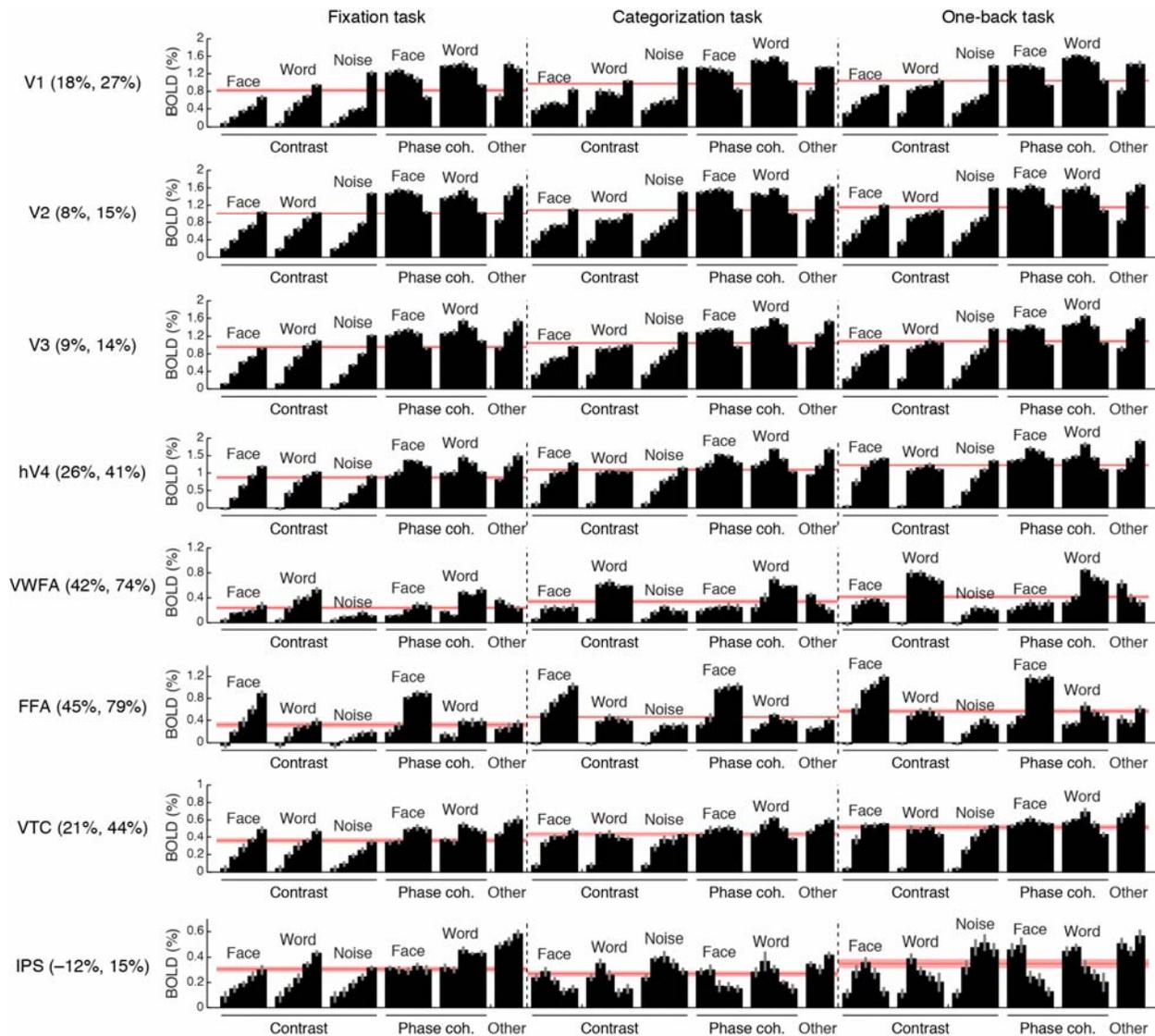
- 1209 Nature Rev Neurosci 3:201–215.
- 1210 Cox DD, Savoy RL (2003) Functional magnetic resonance imaging (fMRI) “brain reading”: detecting and  
1211 classifying distributed patterns of fMRI activity in human visual cortex. *NeuroImage* 19:261–270.
- 1212 Dehaene S, Cohen L (2007) Cultural recycling of cortical maps. *Neuron* 56:384–398.
- 1213 Dehaene S, Cohen L (2011) The unique role of the visual word form area in reading. *Trends in cognitive  
1214 sciences* 15:254–262.
- 1215 Desikan RS, Ségonne F, Fischl B, Quinn BT, Dickerson BC, Blacker D, Buckner RL, Dale AM, Maguire  
1216 RP, Hyman BT, Albert MS, Killiany RJ (2006) An automated labeling system for subdividing the  
1217 human cerebral cortex on MRI scans into gyral based regions of interest. *NeuroImage* 31:968–980.
- 1218 DiCarlo JJ, Zoccolan D, Rust NC (2012) How does the brain solve visual object recognition? *Neuron*  
1219 73:415–434.
- 1220 Donkin C, Brown S, Heathcote A, Wagenmakers E-J (2011) Diffusion versus linear ballistic  
1221 accumulation: different models but the same conclusions about psychological processes? *Psychon  
1222 Bull Rev* 18:61–69.
- 1223 Donner TH, Sagi D, Bonneh YS, Heeger DJ (2013) Retinotopic patterns of correlated fluctuations in  
1224 visual cortex reflect the dynamics of spontaneous perceptual suppression. *J Neurosci* 33:2188–2198.
- 1225 Dumoulin SO, Wandell B (2008) Population receptive field estimates in human visual cortex.  
1226 *NeuroImage* 39:647–660.
- 1227 Feinberg DA, Moeller S, Smith SM, Auerbach E, Ramanna S, Gunther M, Glasser MF, Miller KL,  
1228 Ugurbil K, Yacoub E (2010) Multiplexed echo planar imaging for sub-second whole brain FMRI and  
1229 fast diffusion imaging. Valdes-Sosa PA, ed. *PLoS ONE* 5:e15710.
- 1230 Fischl B (2012) FreeSurfer. *NeuroImage* 62:774–781.
- 1231 Fukushima K (1980) Neocognitron: a self organizing neural network model for a mechanism of pattern  
1232 recognition unaffected by shift in position. *Biological cybernetics* 36:193–202.
- 1233 Gauthier I, Tarr MJ, Anderson AW, Skudlarski P, Gore JC (1999) Activation of the middle fusiform “face  
1234 area” increases with expertise in recognizing novel objects. *Nature neuroscience* 2:568–573.
- 1235 Ghuman AS, Brunet NM, Li Y, Konecky RO, Pyles JA, Walls SA, Destefino V, Wang W, Richardson  
1236 RM (2014) Dynamic encoding of face information in the human fusiform gyrus. *Nat Commun*  
1237 5:5672.
- 1238 Gold JJ, Shadlen MN (2007) The neural basis of decision making. *Annual review of neuroscience*  
1239 30:535–574.
- 1240 Grill-Spector K, Weiner KS (2014) The functional architecture of the ventral temporal cortex and its role  
1241 in categorization. *Nature Rev Neurosci* 15:536–548.
- 1242 Güçlü U, van Gerven MAJ (2015) Deep Neural Networks Reveal a Gradient in the Complexity of Neural  
1243 Representations across the Ventral Stream. *J Neurosci* 35:10005–10014.

- 1244 Hastie T, Tibshirani R, Friedman JH (2001) The elements of statistical learning: data mining, inference,  
1245 and prediction. New York: Springer.
- 1246 Haxby JV, Gobbini MI, Furey ML, Ishai A, Schouten JL, Pietrini P (2001) Distributed and overlapping  
1247 representations of faces and objects in ventral temporal cortex. *Science* 293:2425–2430.
- 1248 Heeger DJ (1992) Normalization of cell responses in cat striate cortex. *Visual neuroscience* 9:181–197.
- 1249 Heeger DJ, Simoncelli EP, Movshon JA (1996) Computational models of cortical visual processing.  
1250 *Proceedings of the National Academy of Sciences of the United States of America* 93:623–627.
- 1251 Heekeren HR, Marrett S, Bandettini P, Ungerleider LG (2004) A general mechanism for perceptual  
1252 decision-making in the human brain. *Nature* 431:859–862.
- 1253 Hubel DH, Wiesel TN (1963) Receptive fields of cells in striate cortex of very young, visually  
1254 inexperienced kittens. *Journal of neurophysiology* 26:994–1002.
- 1255 Huth AG, Nishimoto S, Vu AT, Gallant JL (2012) A continuous semantic space describes the  
1256 representation of thousands of object and action categories across the human brain. *Neuron* 76:1210–  
1257 1224.
- 1258 Itthipuripat S, Ester EF, Deering S, Serences JT (2014) Sensory gain outperforms efficient readout  
1259 mechanisms in predicting attention-related improvements in behavior. *J Neurosci* 34:13384–13398.
- 1260 Jeurissen B, Leemans A, Jones DK, Tournier J-D, Sijbers J (2011) Probabilistic fiber tracking using the  
1261 residual bootstrap with constrained spherical deconvolution. *Hum Brain Mapp* 32:461–479.
- 1262 Jones JP, Palmer LA (1987) The two-dimensional spatial structure of simple receptive fields in cat striate  
1263 cortex. *Journal of neurophysiology* 58:1187–1211.
- 1264 Jozwik KM, Kriegeskorte N, Mur M (2016) Visual features as stepping stones toward semantics:  
1265 Explaining object similarity in IT and perception with non-negative least squares. *Neuropsychologia*  
1266 83:201–226.
- 1267 Kang X, Yund EW, Herron TJ, Woods DL (2007) Improving the resolution of functional brain imaging:  
1268 analyzing functional data in anatomical space. *Magnetic resonance imaging* 25:1070–1078.
- 1269 Kanwisher N (2010) Functional specificity in the human brain: a window into the functional architecture  
1270 of the mind. *Proceedings of the National Academy of Sciences of the United States of America*  
1271 107:11163–11170.
- 1272 Kanwisher N, McDermott J, Chun MM (1997) The fusiform face area: a module in human extrastriate  
1273 cortex specialized for face perception. *J Neurosci* 17:4302–4311.
- 1274 Kanwisher N, Wojciulik E (2000) Visual attention: insights from brain imaging. *Nature Rev Neurosci*  
1275 1:91–100.
- 1276 Kay KN, Naselaris T, Prenger RJ, Gallant JL (2008) Identifying natural images from human brain  
1277 activity. *Nature* 452:352–355.
- 1278 Kay KN, Rokem A, Winawer J, Dougherty RF, Wandell B (2013a) GLMdenoise: a fast, automated  
1279 technique for denoising task-based fMRI data. *Front Neurosci* 7:247.

- 1280 Kay KN, Winawer J, Mezer A, Wandell B (2013b) Compressive spatial summation in human visual  
1281 cortex. *Journal of neurophysiology* 110:481–494.
- 1282 Kay KN, Winawer J, Rokem A, Mezer A, Wandell B (2013c) A two-stage cascade model of BOLD  
1283 responses in human visual cortex. Diedrichsen J, ed. *PLoS computational biology* 9:e1003079.
- 1284 Kayser AS, Buchsbaum BR, Erickson DT, D'Esposito M (2010a) The functional anatomy of a perceptual  
1285 decision in the human brain. *Journal of neurophysiology* 103:1179–1194.
- 1286 Kayser AS, Erickson DT, Buchsbaum BR, D'Esposito M (2010b) Neural representations of relevant and  
1287 irrelevant features in perceptual decision making. *J Neurosci* 30:15778–15789.
- 1288 Khaligh-Razavi S-M, Kriegeskorte N (2014) Deep supervised, but not unsupervised, models may explain  
1289 IT cortical representation. Diedrichsen J, ed. *PLoS computational biology* 10:e1003915.
- 1290 Kriegeskorte N, Mur M, Ruff DA, Kiani R, Bodurka J, Esteky H, Tanaka K, Bandettini PA (2008)  
1291 Matching categorical object representations in inferior temporal cortex of man and monkey. *Neuron*  
1292 60:1126–1141.
- 1293 Kriegeskorte N, Simmons WK, Bellgowan PSF, Baker CI (2009) Circular analysis in systems  
1294 neuroscience: the dangers of double dipping. *Nature neuroscience* 12:535–540.
- 1295 Lauritzen TZ, D'Esposito M, Heeger DJ, Silver MA (2009) Top-down flow of visual spatial attention  
1296 signals from parietal to occipital cortex. *Journal of vision* 9:18.1–14.
- 1297 Luck SJ, Chelazzi L, Hillyard SA, Desimone R (1997) Neural mechanisms of spatial selective attention in  
1298 areas V1, V2, and V4 of macaque visual cortex. *Journal of neurophysiology* 77:24–42.
- 1299 McAdams CJ, Maunsell JH (1999) Effects of attention on orientation-tuning functions of single neurons  
1300 in macaque cortical area V4. *J Neurosci* 19:431–441.
- 1301 Moeller S, Yacoub E, Olman CA, Auerbach E, Strupp J, Harel N, Ugurbil K (2010) Multiband multislice  
1302 GE-EPI at 7 tesla, with 16-fold acceleration using partial parallel imaging with application to high  
1303 spatial and temporal whole-brain fMRI. *Magn Reson Med* 63:1144–1153.
- 1304 Murray SO, He S (2006) Contrast invariance in the human lateral occipital complex depends on attention.  
1305 *Curr Biol* 16:606–611.
- 1306 Nasr S, Echavarria CE, Tootell RBH (2014) Thinking outside the box: rectilinear shapes selectively  
1307 activate scene-selective cortex. *J Neurosci* 34:6721–6735.
- 1308 O'Reilly JX, Woolrich MW, Behrens TEJ, Smith SM, Johansen-Berg H (2012) Tools of the trade:  
1309 psychophysiological interactions and functional connectivity. *Social cognitive and affective*  
1310 *neuroscience* 7:604–609.
- 1311 Ohayon S, Freiwald WA, Tsao DY (2012) What makes a cell face selective? The importance of contrast.  
1312 *Neuron* 74:567–581.
- 1313 Pelli DG (1997) The VideoToolbox software for visual psychophysics: transforming numbers into  
1314 movies. *Spat Vis* 10:437–442.
- 1315 Pitzalis S, Fattori P, Galletti C (2012) The functional role of the medial motion area V6. *Front Behav*

- 1316 Neurosci 6:91.
- 1317 Price CJ, Devlin JT (2011) The interactive account of ventral occipitotemporal contributions to reading.  
1318 Trends in cognitive sciences 15:246–253.
- 1319 Rainer G, Augath M, Trinath T, Logothetis NK (2001) Nonmonotonic noise tuning of BOLD fMRI signal  
1320 to natural images in the visual cortex of the anesthetized monkey. Curr Biol 11:846–854.
- 1321 Ratcliff R (1978) A theory of memory retrieval. Psychological review 85:59.
- 1322 Reich L, Szwed M, Cohen L, Amedi A (2011) A ventral visual stream reading center independent of  
1323 visual experience. Curr Biol 21:363–368.
- 1324 Ress D, Backus BT, Heeger DJ (2000) Activity in primary visual cortex predicts performance in a visual  
1325 detection task. Nature neuroscience 3:940–945.
- 1326 Ress D, Heeger DJ (2003) Neuronal correlates of perception in early visual cortex. Nature neuroscience  
1327 6:414–420.
- 1328 Reynolds JH, Heeger DJ (2009) The normalization model of attention. Neuron 61:168–185.
- 1329 Reynolds JH, Pasternak T, Desimone R (2000) Attention increases sensitivity of V4 neurons. Neuron  
1330 26:703–714.
- 1331 Rissman J, Gazzaley A, D'Esposito M (2004) Measuring functional connectivity during distinct stages of  
1332 a cognitive task. NeuroImage 23:752–763.
- 1333 Rolls ET (2012) Invariant Visual Object and Face Recognition: Neural and Computational Bases, and a  
1334 Model, VisNet. Front Comput Neurosci 6:35.
- 1335 Saalmann YB, Pigarev IN, Vidyasagar TR (2007) Neural mechanisms of visual attention: how top-down  
1336 feedback highlights relevant locations. Science 316:1612–1615.
- 1337 Schira MM, Tyler CW, Breakspear M, Spehar B (2009) The foveal confluence in human visual cortex. J  
1338 Neurosci 29:9050–9058.
- 1339 Sereno AB, Maunsell JH (1998) Shape selectivity in primate lateral intraparietal cortex. Nature 395:500–  
1340 503.
- 1341 Serre T, Kreiman G, Kouh M, Cadieu C, Knoblich U, Poggio T (2007) A quantitative theory of  
1342 immediate visual recognition. Progress in brain research 165:33–56.
- 1343 Shadlen MN, Newsome WT (2001) Neural basis of a perceptual decision in the parietal cortex (area LIP)  
1344 of the rhesus monkey. Journal of neurophysiology 86:1916–1936.
- 1345 Striem-Amit E, Cohen L, Dehaene S, Amedi A (2012) Reading with sounds: sensory substitution  
1346 selectively activates the visual word form area in the blind. Neuron 76:640–652.
- 1347 Takemura H, Rokem A, Winawer J, Yeatman JD, Wandell B, Pestilli F (2016) A Major Human White  
1348 Matter Pathway Between Dorsal and Ventral Visual Cortex. Cereb Cortex 26:2205–2214.
- 1349 Tournier J-D, Calamante F, Connelly A (2007) Robust determination of the fibre orientation distribution

- 1350 in diffusion MRI: non-negativity constrained super-resolved spherical deconvolution. *NeuroImage*  
1351 35:1459–1472.
- 1352 Twomey T, Kawabata Duncan KJ, Price CJ, Devlin JT (2011) Top-down modulation of ventral occipito-  
1353 temporal responses during visual word recognition. *NeuroImage* 55:1242–1251.
- 1354 Vinckier F, Dehaene S, Jobert A, Dubus JP, Sigman M, Cohen L (2007) Hierarchical coding of letter  
1355 strings in the ventral stream: dissecting the inner organization of the visual word-form system.  
1356 *Neuron* 55:143–156.
- 1357 Wandell B (1999) Computational neuroimaging of human visual cortex. *Annual review of neuroscience*  
1358 22:145–173.
- 1359 Wandell B, Rauschecker AM, Yeatman JD (2012) Learning to see words. *Annual review of psychology*  
1360 63:31–53.
- 1361 Wang L, Mruczek REB, Arcaro MJ, Kastner S (2014) Probabilistic Maps of Visual Topography in  
1362 Human Cortex. *Cereb Cortex*:bhu277.
- 1363 Weiner KS, Golarai G, Caspers J, Chuapoco MR, Mohlberg H, Zilles K, Amunts K, Grill-Spector K  
1364 (2014) The mid-fusiform sulcus: a landmark identifying both cytoarchitectonic and functional  
1365 divisions of human ventral temporal cortex. *NeuroImage* 84:453–465.
- 1366 Weiner KS, Grill-Spector K (2010) Sparsely-distributed organization of face and limb activations in  
1367 human ventral temporal cortex. *NeuroImage* 52:1559–1573.
- 1368 Weiner KS, Grill-Spector K (2011) Not one extrastriate body area: using anatomical landmarks, hMT+,  
1369 and visual field maps to parcellate limb-selective activations in human lateral occipitotemporal  
1370 cortex. *NeuroImage* 56:2183–2199.
- 1371 Wu MC, David SV, Gallant JL (2006) Complete functional characterization of sensory neurons by system  
1372 identification. *Annual review of neuroscience* 29:477–505.
- 1373 Yamins DLK, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ (2014) Performance-optimized  
1374 hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National*  
1375 *Academy of Sciences of the United States of America* 111:8619–8624.
- 1376 Yeatman JD, Rauschecker AM, Wandell B (2013) Anatomy of the visual word form area: adjacent  
1377 cortical circuits and long-range white matter connections. *Brain Lang* 125:146–155.
- 1378 Yeatman JD, Weiner KS, Pestilli F, Rokem A, Mezer A, Wandell B (2014) The vertical occipital  
1379 fasciculus: a century of controversy resolved by in vivo measurements. *Proceedings of the National*  
1380 *Academy of Sciences of the United States of America* 111:E5214–E5223.
- 1381 Yue X, Cassidy BS, Devaney KJ, Holt DJ, Tootell RBH (2011) Lower-level stimulus features strongly  
1382 influence responses in the fusiform face area. *Cereb Cortex* 21:35–47.
- 1383 Zetzsche C, Krieger G, Wegmann B (1999) The atoms of vision: Cartesian or polar? *J Opt Soc Am A,*  
1384 *JOSAA* 16:1554–1565.
- 1385



1386

1387

1388 **Figure 1–figure supplement 1. Comprehensive summary of fMRI measurements.** Black bars indicate

1389 responses (beta weights) evoked by different stimuli and tasks. Red lines indicate the average response

1390 across stimuli, computed separately for each task. Error bars indicate bootstrapped 68% CIs (resampling

1391 subjects with replacement). Percentages in ROI labels indicate the strength of the response observed

1392 during the categorization and one-back tasks relative to the fixation task. For example, in FFA, the

1393 average response across stimuli during the one-back task is 79% stronger than the average response

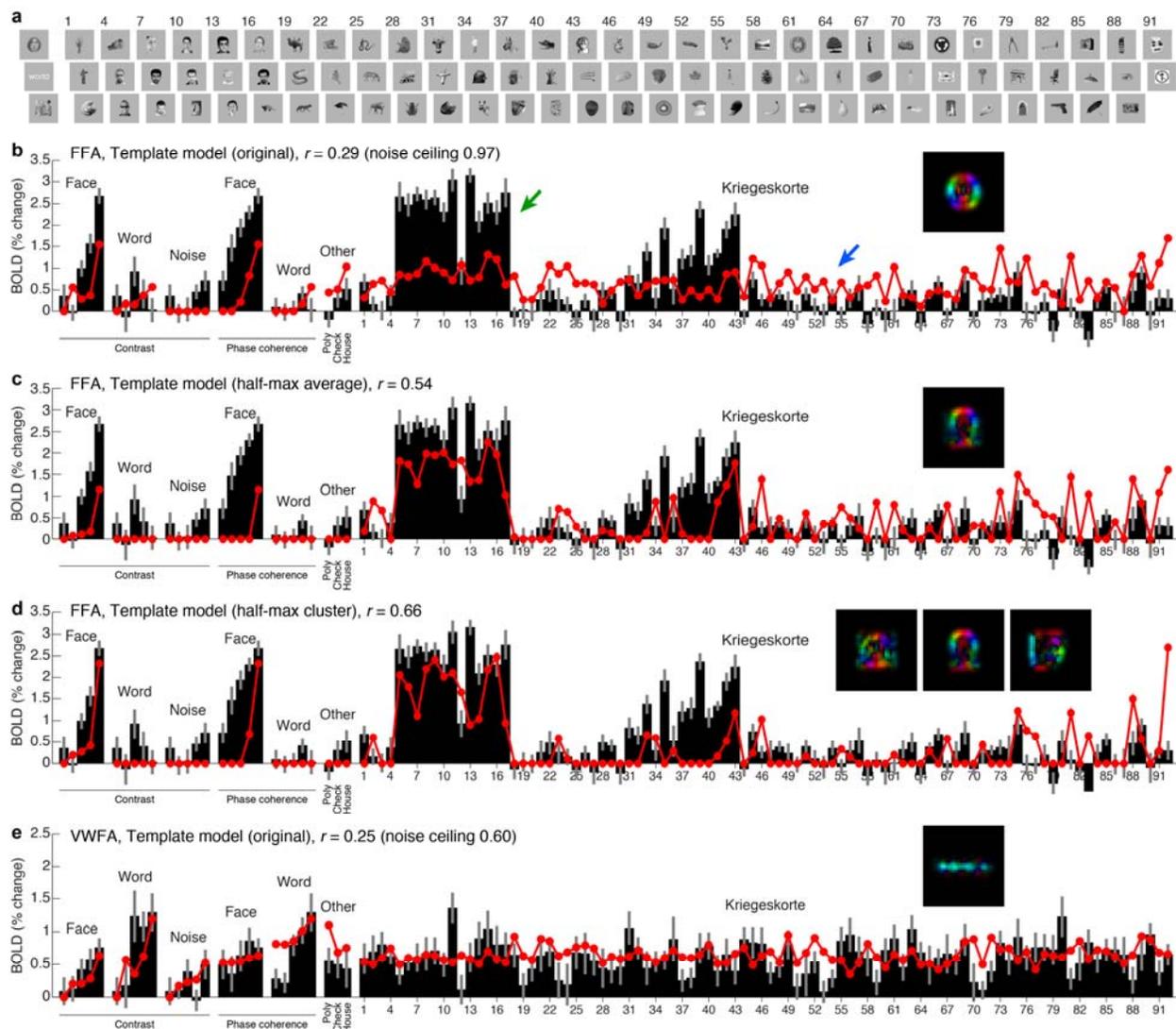
1394 across stimuli during the fixation task. Task effects are substantially stronger in VWFA and FFA than in

1395 early visual areas V1–V3. The larger apparent task modulation in V1 compared to V2 and V3 might due

1396 to small eye movements that may have been made during the categorization and one-back tasks. Our

1397 interpretation of the observed IPS activity during the fixation task is that this activity reflects the decision-

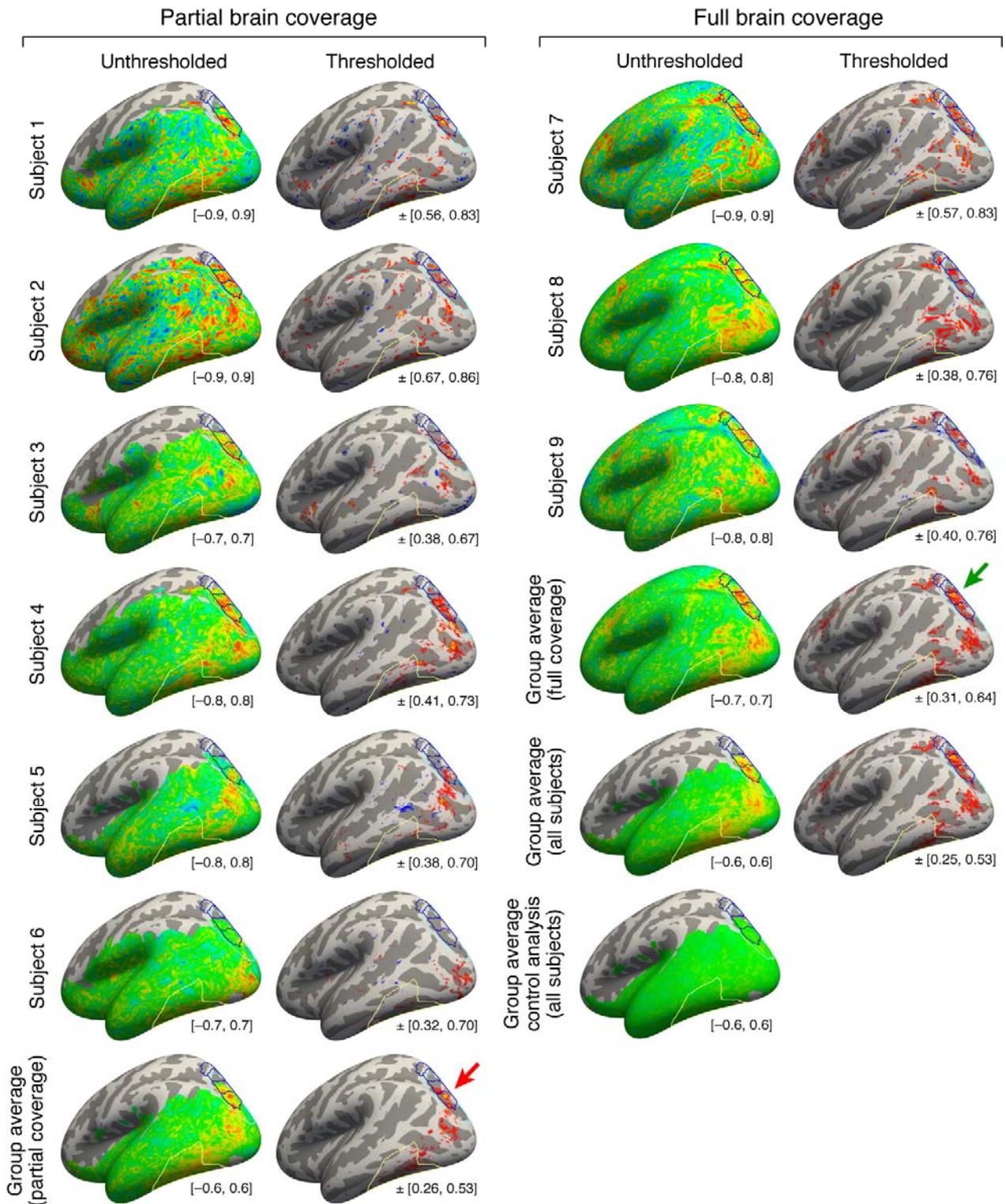
1398 making process involved in judging the color of the fixation dot. Support for this interpretation comes  
1399 from the fact that the root-mean-square contrast of the stimuli, computed over a small region surrounding  
1400 the fixation dot ( $0.36^\circ \times 0.36^\circ$ ), correlates strongly with IPS responses during the fixation task ( $r = 0.86$ ).  
1401  
1402



1403  
1404

1405 **Figure 2–figure supplement 1. Testing the Template model on a wide range of stimuli.** (a) *Stimuli.*  
 1406 We collected an additional dataset consisting of 92 images from a previous study by Kriegeskorte et al.  
 1407 (Kriegeskorte et al., 2008) (all images shown), along with 22 images from the original experiment (three  
 1408 images shown). We assessed model accuracy using 20-fold cross-validation across stimuli (see Methods  
 1409 for details). (b) *Performance of Template model (original).* Black bars indicate data from FFA, with error  
 1410 bars indicating 68% CIs (error across trials). Red lines and red dots indicate model predictions. Inset  
 1411 shows the category template used in the model. The model performs poorly. (c) *Performance of Template*  
 1412 *model (half-max average).* This model derives the category template by computing (in the V1-like  
 1413 representation) the centroid of all stimuli in the training set that evoke at least half of the maximum  
 1414 response. Performance improves. (d) *Performance of Template model (half-max cluster).* This model  
 1415 derives multiple category templates by performing *k*-means clustering (in the V1-like representation) on

1416 all stimuli in the training set that evoke at least half of the maximum response. Performance further  
1417 improves, resolving both underprediction of responses (e.g. green arrow in panel b) and overprediction of  
1418 responses (e.g. blue arrow in panel b). (e) *Results for VWFA*. Similar responses are observed across the 92  
1419 Kriegeskorte images. Responses are well predicted by the original Template model, up to the level of  
1420 measurement noise in this region.  
1421



1422

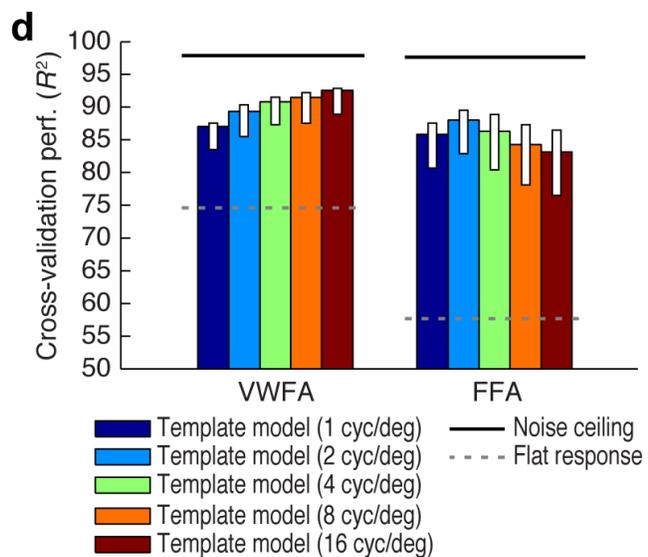
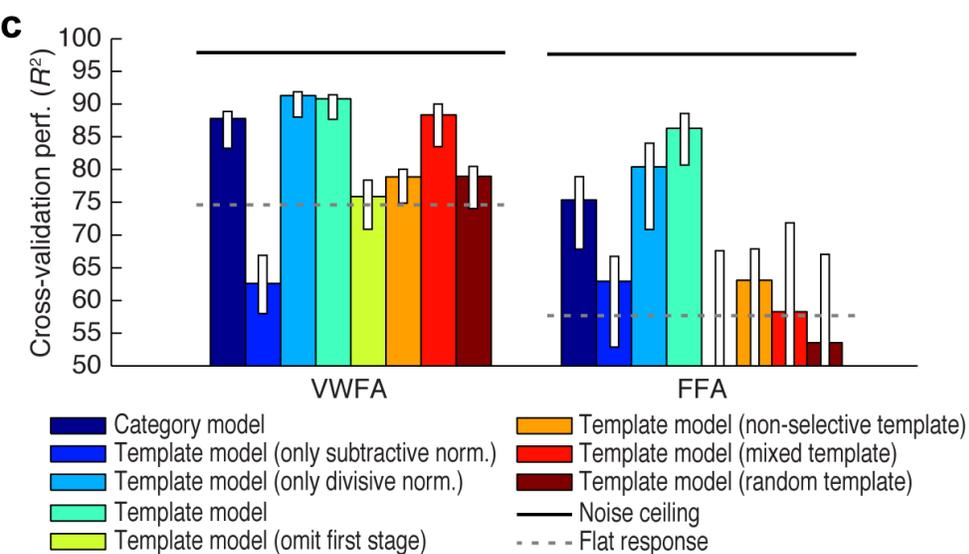
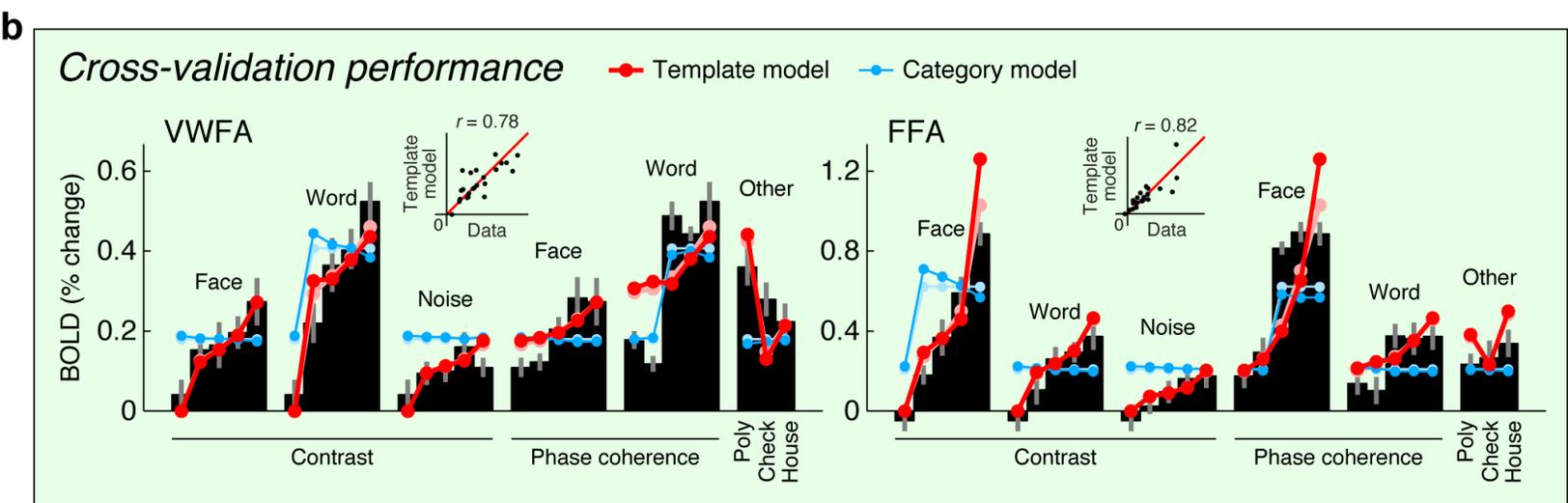
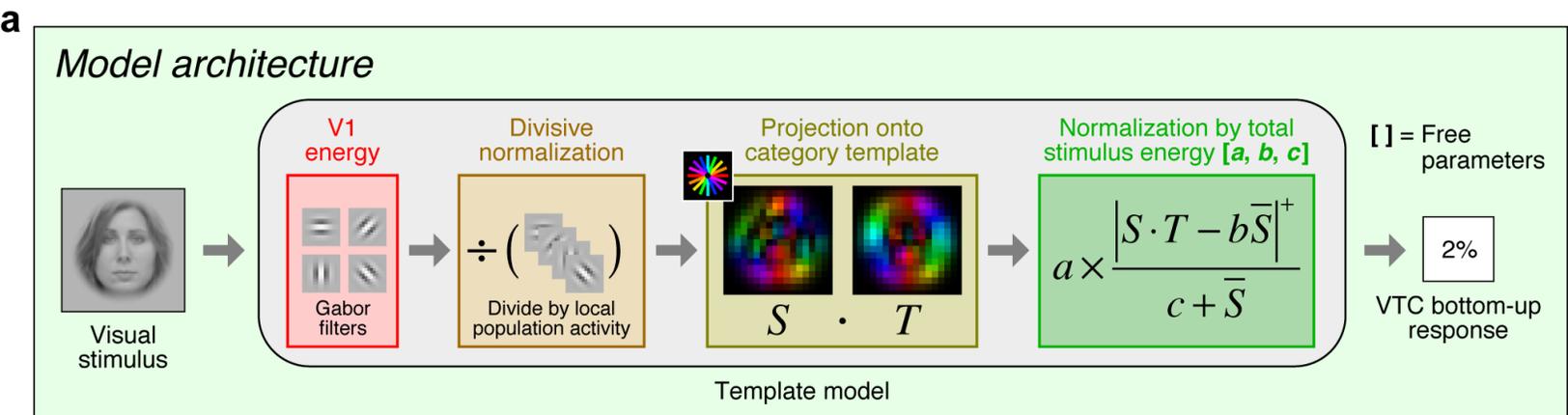
1423

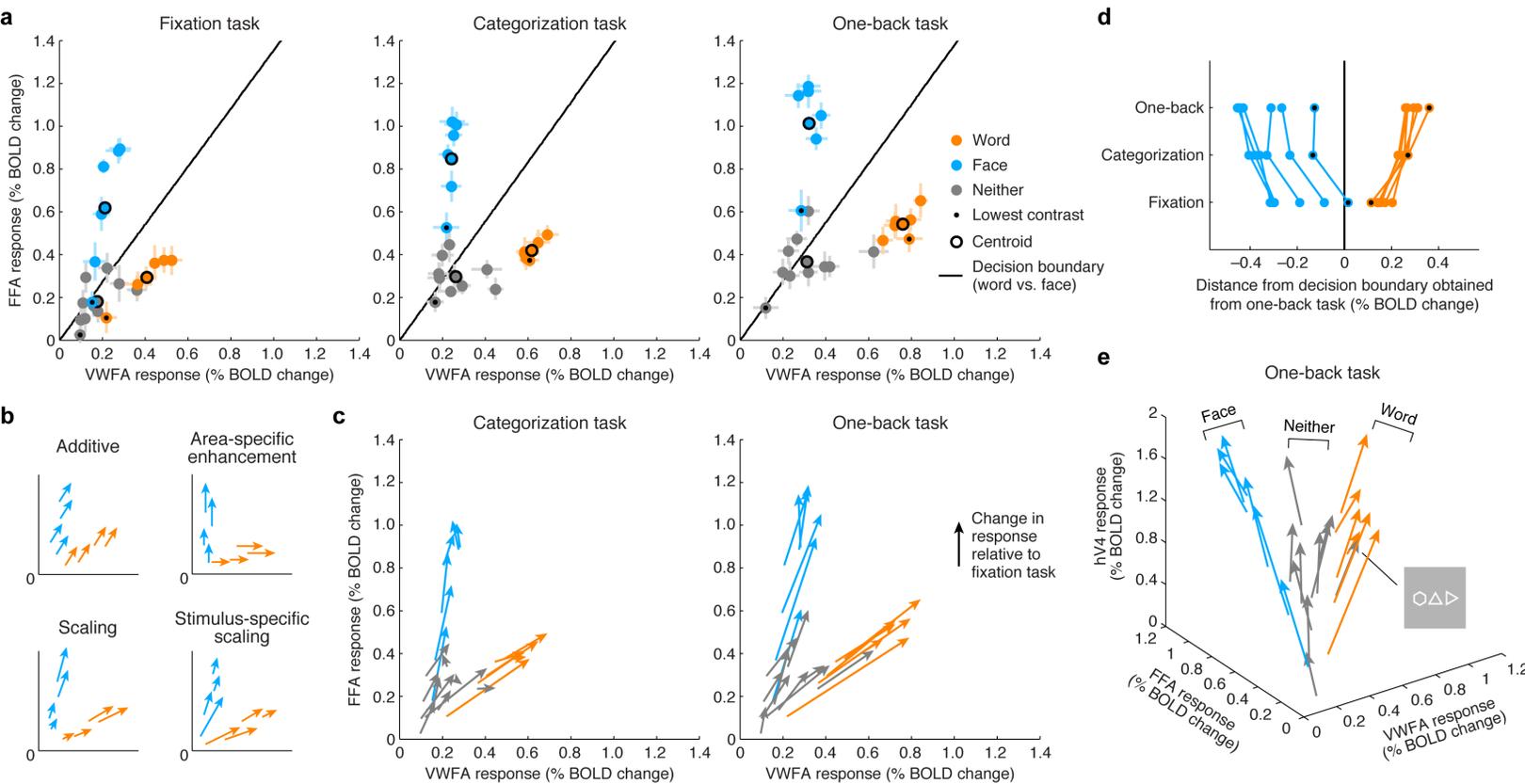
1424 **Figure 4–figure supplement 1. Maps of top-down connectivity to VTC.** This figure shows thresholded and unthresholded maps for individual subjects and group averages (same format as **Figure 4b**; all maps

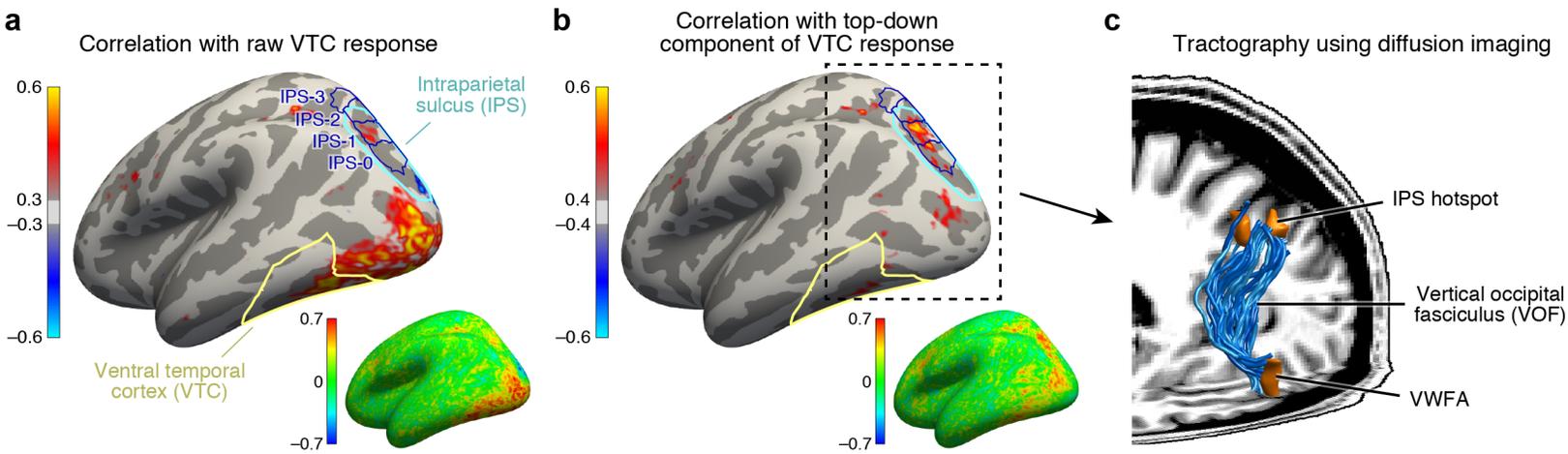
1425

1426 shown on the *fsaverage* surface). At the lower right of each map is the range of values used for the  
1427 colormap. The left two columns show results obtained for the six subjects with partial brain coverage.  
1428 Group average results for these subjects are shown in the last row. The right two columns show results  
1429 obtained for the three subjects with full brain coverage. Group average results for these subjects are  
1430 shown in the third to last row. Group average results for all subjects are shown in the second to last row.  
1431 The last row shows results obtained from a control analysis in which we generate individual-subject maps  
1432 by correlating cortical responses with random Gaussian noise and then average these maps across  
1433 subjects. This control analysis produces no substantial correlations. Notice that the peak correlation is  
1434 found in and around IPS-0/1 for both the group of subjects with partial brain coverage (red arrow) and the  
1435 group of subjects with full brain coverage (green arrow). Some variability in the location of the peak  
1436 correlation is expected given that there are limits on the degree to which functional areas can be aligned  
1437 across subjects based solely on anatomical features.  
1438

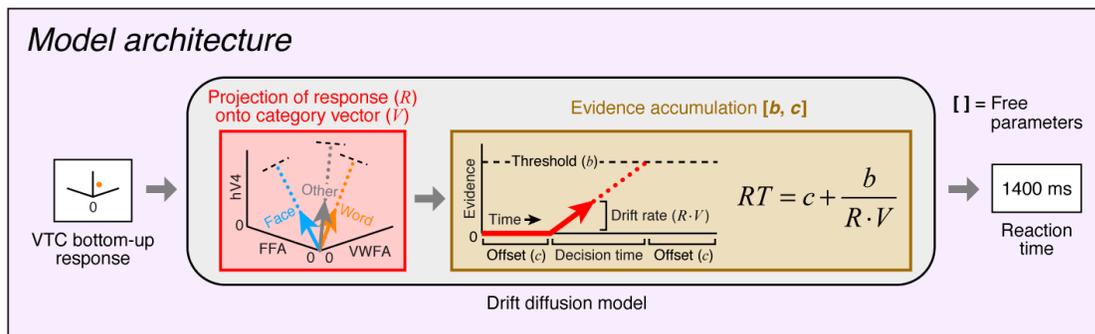
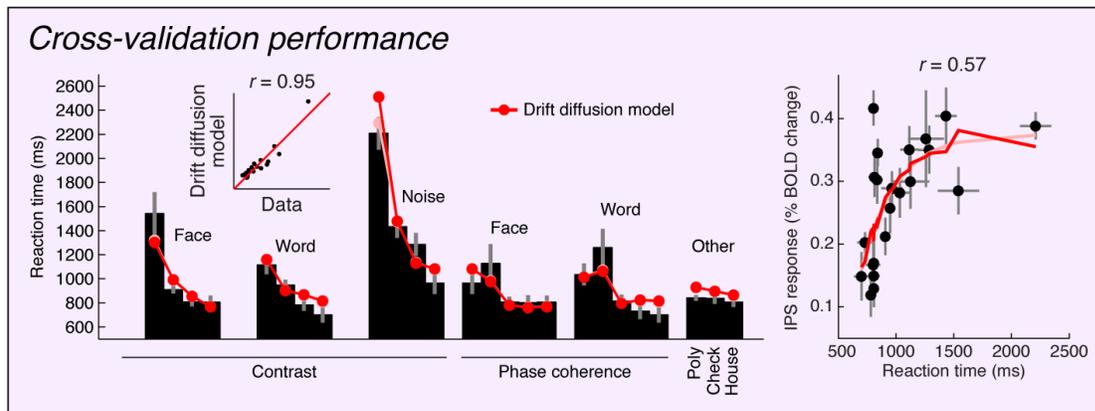
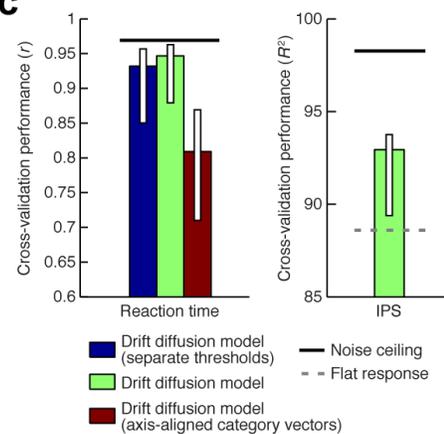
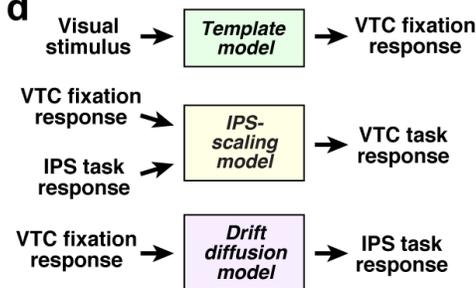


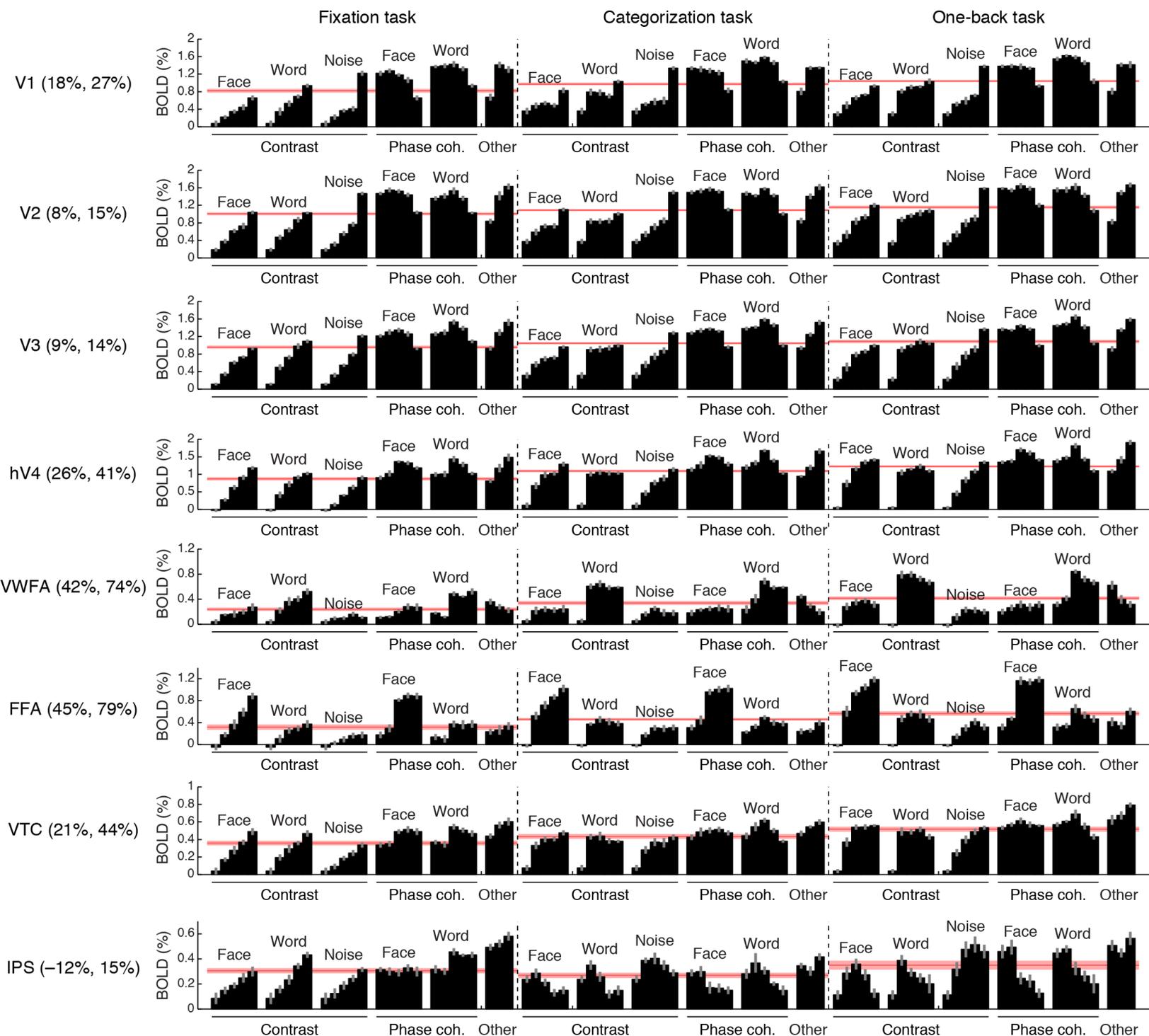


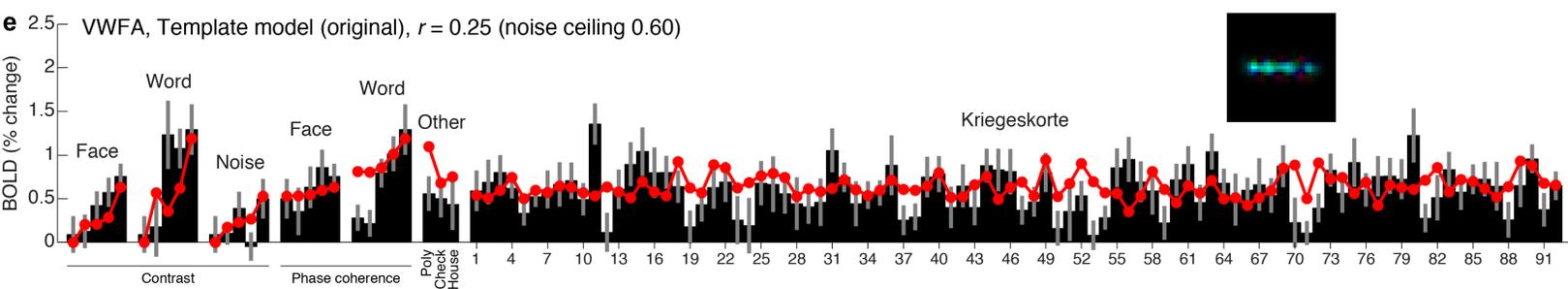
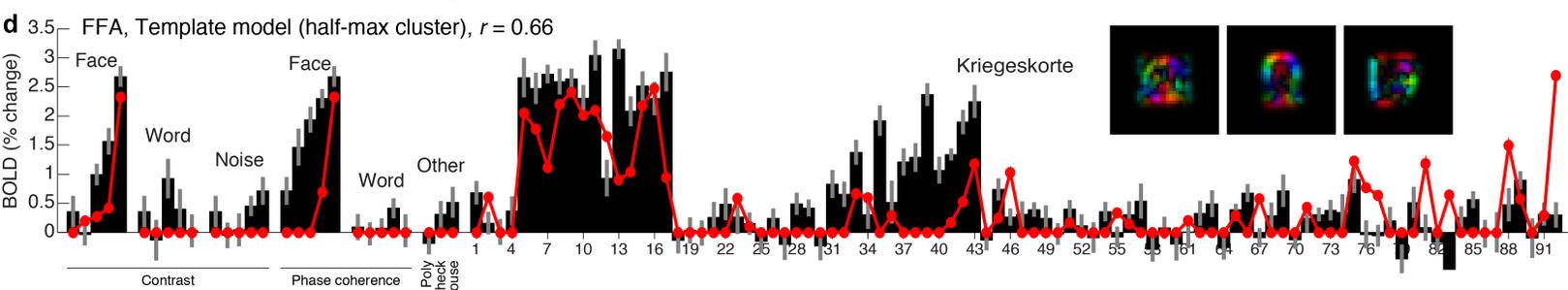
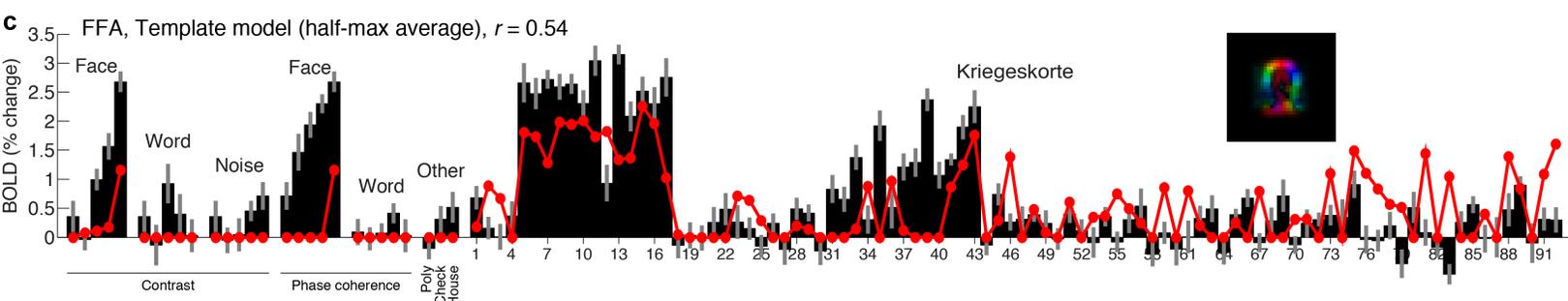
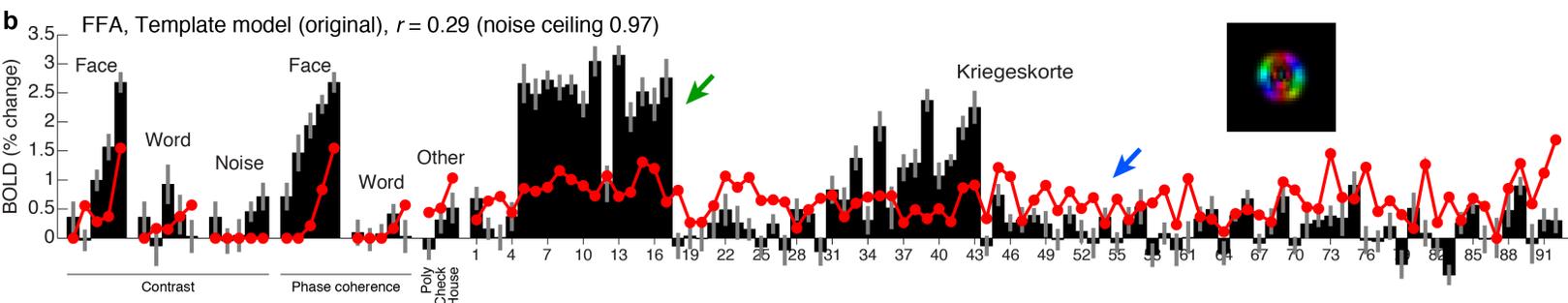
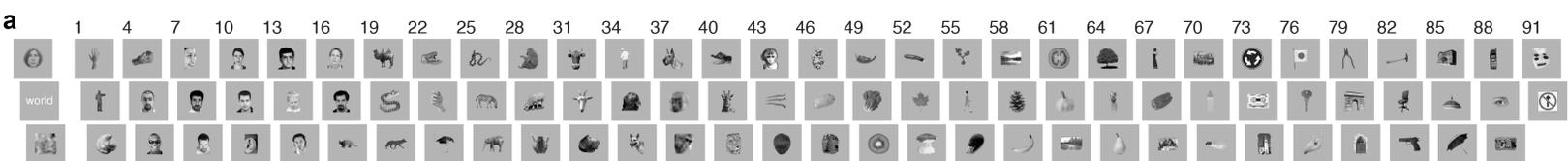




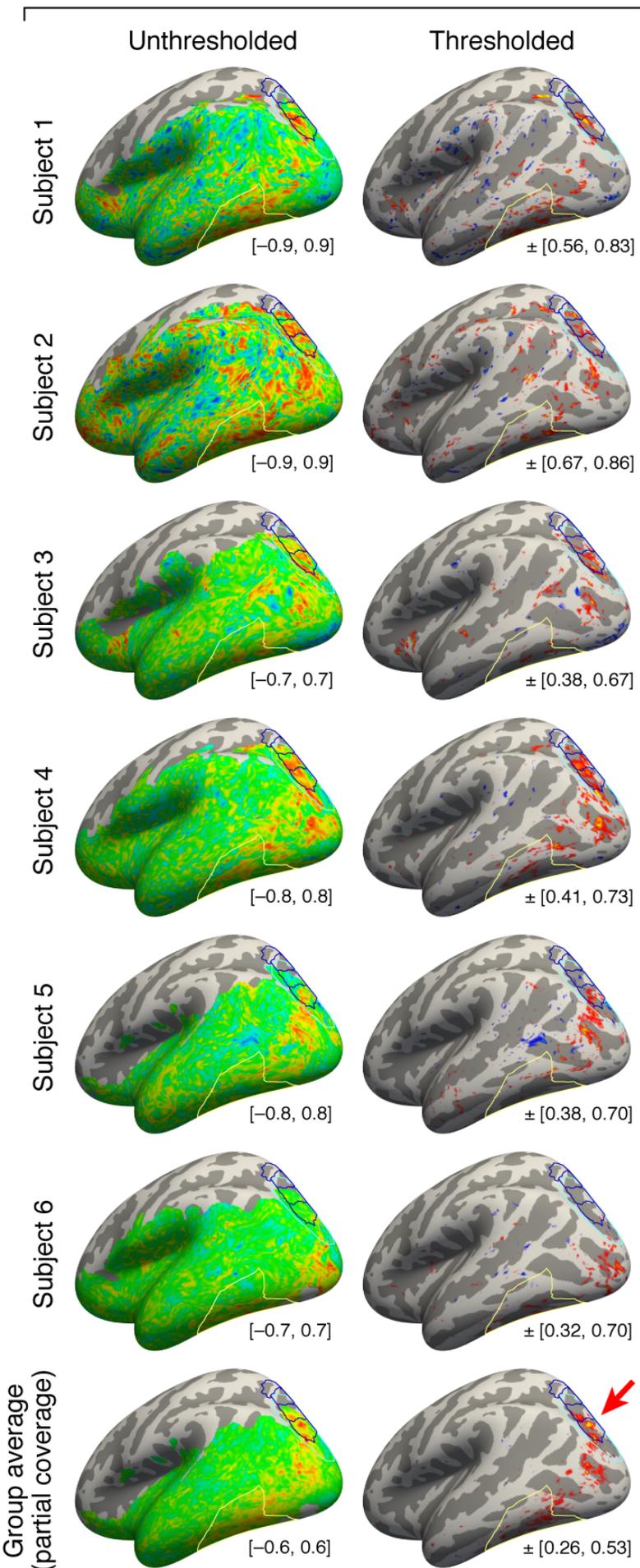


**a****b****c****d**





## Partial brain coverage



## Full brain coverage

